

Posterior contraction of the population polytope in finite admixture models ¹

XuanLong Nguyen
xuanlong@umich.edu

Technical report 528
Department of Statistics
University of Michigan

First version, May 24, 2012
Revised, Dec 4, 2012

Abstract

We study the posterior contraction behavior of the latent population structure that arises in admixture models as the amount of data increases. An admixture model — alternatively known as a topic model — specifies k populations, each of which is characterized by a Δ^d -valued vector of frequencies for generating a set of discrete values in $\{0, 1, \dots, d\}$. The population polytope is defined as the convex hull of the k frequency vectors. Under the admixture specification, each of m individuals generates an i.i.d. frequency vector according to a probability distribution defined on the (unknown) population polytope G_0 , and then generates n data points according to the sampled frequency vector. Rates of posterior contraction are established with respect to Hausdorff metric and a minimum matching Euclidean metric defined on population polytopes, as the amount of data $m \times n$ tends to infinity. Minimax lower bounds are also established. Tools developed include posterior asymptotics of hierarchical models with $m \times n$ data, and arguments from convex geometry.

1 Introduction

We study a class of hierarchical mixture models for categorical data known as the admixtures, which were independently developed in the landmark papers by Pritchard, Stephens and Donnelly [Pritchard et al., 2000] and Blei, Ng and Jordan [Blei et al., 2003]. The former set of authors applied their modeling to population genetics, while the latter considered applications in text processing and computer vision, where their models are more widely

¹ AMS 2000 subject classification. Primary 62F15, 62G05; secondary 62G20.

Key words and phrases: latent mixing measures, convex polytope, population structure, topic simplex, Bayesian estimation, posterior consistency, rates of convergence, latent Dirichlet allocation, Hausdorff metric, convex geometry

This research was supported in part by NSF grants CCF-1115769 and OCI-1047871. The author thank Qiaozhu Mei and Jian Tang for stimulating discussions on topic models, which helped to motivate this work. Thanks also go to helpful comments by anonymous referees on an earlier draft.

known as the *latent Dirichlet allocation* model, or a topic model. Admixture modeling has been applied to and extended in a vast number of fields of engineering and sciences — in fact, the Google scholar pages for these two original papers alone combine for more than a dozen thousands of citations. In spite of their wide uses, asymptotic behavior of hierarchical models such as the admixtures remains largely unexplored, to the best of our knowledge.

A finite admixture model posits that there are k populations, each of which is characterized by a Δ^d -valued vector θ_j of frequencies for generating a set of discrete values $\{0, 1, \dots, d\}$, for $j = 1, \dots, k$. Here, Δ^d is the d -dimensional probability simplex. A sampled individual may have mixed ancestry and as a result inherits some fraction of its values from each of its ancestral populations. Thus, an individual is associated with a proportion vector $\beta = (\beta_1, \dots, \beta_k) \in \Delta^{k-1}$, where β_j denotes the proportion of the individual's data that are generated according to population j 's frequency vector θ_j . This yields a vector of frequencies $\eta = \sum_{j=1}^k \beta_j \theta_j \in \Delta^d$ associated with that individual. In most applications, one does not observe η directly, but rather an i.i.d. sample generated from a multinomial distribution parameterized by η . The collection of $\theta_1, \dots, \theta_k$ is referred to as the *population structure* in the admixture. In population genetics modeling, θ_j represents the allele frequencies at each locus in an individual's genome from the j -th population. In text document modeling, θ_j represents the frequencies of words generated by the j -th topic, while an individual is a document, i.e., a collection of words. In computer vision, θ_j represents the frequencies of objects generated by the j -th scenery topic, while an individual is a natural image, i.e., a collection of scenery objects. The primary interest is the inference of the population structure on the basis of sampled data. In a Bayesian estimation setting, the population structure is assumed random and endowed with a prior distribution — accordingly one is interested in the behavior of the posterior distribution of the population structure given the available data.

The goal of this paper is to obtain contraction rates of the posterior distribution of the latent population structure that arises in admixture models, as the amount of data increases. Admixture models present a canonical mixture model for categorical data in which the population structure provides the support for the mixing measure. Existing works on convergence behavior of mixing measures in a mixture model are quite rare, in either frequentist or Bayesian estimation literature. Chen provided the optimal convergence rate of mixing measures in several finite mixtures for univariate data Chen [1995] (see also Ishwaran et al. [2001]). Recent progress on multivariate mixture models include papers by Rousseau and Mengersen Rousseau and Mengersen [2011] and Nguyen Nguyen [2012]. In Nguyen [2012] posterior contraction rates of mixing measures in several finite and infinite mixture models for multivariate and continuous data were obtained. Toussile and Gassiat established consistency of a penalized MLE procedure for a finite admixture model Toussile and Gassiat [2009]. This issue has also attracted increased attention in machine learning. Recent papers by Arora et al Arora et al. [2012] and Anandkumar et al Anandkumar et al. [2012] study convergence properties of certain computationally efficient learning algorithms based on matrix factorization techniques.

There are a number of questions that arise in the convergence analysis of admixture

models for categorical data. The first question is to find a suitable metric in order to establish rates of convergence. It would be ideal to establish convergence for each individual element θ_i , for $i = 1, \dots, k$. This is a challenging due to the problems of identifiability. A (relatively) minor issue is known as “label-switching” problem. That is, one can identify the collection of θ_i ’s only up to a permutation. A deeper problem is that any θ_j that can be expressed as a convex combination of the others $\theta_{j'}$ for $j' \neq j$ may be difficult to identify, estimate, and analyze. To get around this difficulty, we propose to study the convergence of population structure variables through its convex hull $G = \text{conv}(\theta_1, \dots, \theta_k)$, which shall be referred to as the *population polytope*. Convergence of convex polytopes can be evaluated in terms of Hausdorff metric $d_{\mathcal{H}}$, a metric commonly utilized in convex geometry Schneider [1993]. Moreover, under some geometric identifiability conditions, it can be shown that convergence in Hausdorff metric entails convergence of all extreme points of the polytope via a minimum-matching distance metric (defined in Section 2). This is the theory we aim for in this paper. Convergence behavior of (the posterior of) non-extreme points among $\theta_1, \dots, \theta_k$ remains unresolved as of this writing.

The second question in an asymptotic study of a hierarchical model is how to address multiple quantities that define the amount of empirical data. The admixture model we consider has two asymptotic quantities that play asymmetric roles — m is the number of individuals, and n is the number of data points associated with each individual. Both m and n are allowed to increase to infinity. A simple way to think about this asymptotic setting is to let m go to infinity, while $n := n(m)$ tends to infinity at a certain rate which may be constrained with respect to m . Let Π be a prior distribution on variables $\theta_1, \dots, \theta_k$. The goal is to derive a vanishing sequence of $\delta_{m,n}$, depending on both m and n , such that the posterior distribution of the θ_i ’s satisfies, for some sufficiently large constant C ,

$$\Pi\left(d_{\mathcal{H}}(G, G_0) \geq C\delta_{m,n} \middle| \mathcal{S}_{[n]}^{[m]}\right) \rightarrow 0$$

in $P_{\mathcal{S}_{[n]}|G_0}^m$ -probability as $m \rightarrow \infty$ and $n = n(m) \rightarrow \infty$ suitably. Here, $P_{\mathcal{S}_{[n]}|G_0}^m$ denotes the true distribution associated with population polytope G_0 that generates a given $m \times n$ data set $\mathcal{S}_{[n]}^{[m]}$. As mentioned, $\delta_{m,n}$ is also the posterior contraction rate for the extreme points among population structure variables $\theta_1, \dots, \theta_k$.

Overview of results. Suppose that $n \rightarrow \infty$ at a rate constrained by $\log m < n$ and $\log n = o(m)$. In an overfitted setting, i.e., when the true population polytope may have less than k extreme points, we show that under some mild identifiability conditions the posterior contraction rate in either Hausdorff or minimum-matching distance metric is

$\delta_{m,n} \asymp \left[\frac{\log m}{m} + \frac{\log n}{n} + \frac{\log n}{m} \right]^{\frac{1}{2(p+\alpha)}}$, where $p = (k-1) \wedge d$ is the intrinsic dimension of the population polytope while α denotes the regularity level near boundary of the support of the density function for η . On the other hand, if either the true population polytope is known to have exactly k extreme points, or if the pairwise distances among the extreme points are bounded from below by a known positive constant, then the contraction rate is

improved to a parametric rate $\delta_{m,n} \asymp \left[\frac{\log m}{m} + \frac{\log n}{n} + \frac{\log n}{m} \right]^{\frac{1}{2(1+\alpha)}}$.

The constraints on $n = n(m)$, and the appearance of quantity $\log n/m$ in the convergence rate are quite interesting. Both the constraints and the derived rate are rooted in a condition on the required thickness of the prior support of the marginal densities of the data and an upper bound on the entropy of the space of such densities. This suggests an interesting interaction between layers in the latent hierarchy of the admixture model worthy of further investigation. For instance, it is not clear whether posterior consistency continues to hold if n falls outside of the specified range, and what effects this has on convergence rates, with or without additional assumptions on the data. This appears quite difficult with our present set of techniques.

We also establish minimax lower bounds for both settings. In the overfitted setting, the obtained lower bound is $(mn)^{-1/(q+\alpha)}$, where $q = \lfloor k/2 \rfloor \wedge d$, unless additional constraints are imposed on the prior. Although this lower bound does not quite match with the posterior contraction rate, the two are qualitatively comparable and both notably dependent on dimensionality d . In particular, if $n \asymp m$, and $k \geq 2d$, the posterior contraction rate becomes $(\log m/m)^{-\frac{1}{2(d+\alpha)}}$. Compare this to the lower bound $m^{-2/(d+\alpha)}$, whose exponent differs by only a factor of 4.

Method of proofs and tools. The general framework of posterior asymptotics for density estimation has been well-established Ghosal et al. [2000], Shen and Wasserman [2001] (see also Barron et al. [1999], Ghosh and Ramamoorthi [2002], Walker [2004], Ghosal and van der Vaart [2007], Walker et al. [2007]). This framework continues to be very useful, but the analysis of mixing measure estimation in multi-level models presents distinct new challenges. In Section 4 we shall formulate an abstract theorem (Theorem 4) on posterior contraction of latent variables of interest in an admixture model, given $m \times n$ data, by reposing on the framework of Ghosal et al. [2000] (see also Nguyen [2012]). The main novelty here is that we work on the space of latent variables (e.g., space of latent population structures endowed with Hausdorff or comparable metric) as opposed to the space of data densities endowed with Hellinger metric. A basic quantity is the *Hellinger information* of the Hausdorff metric for a given subset of polytopes. Indeed, the Hellinger information is a fundamental quantity running through the analysis, which ties together the amount of data m and n — key quantities that are associated with different levels in the model hierarchy.

The bulk of the paper is devoted to establishing properties of the Hellinger information, which are fed into Theorem 4 so as to obtain concrete convergence rates. This is achieved through a number inequalities which illuminate the relationship between Hausdorff distance of a given pair of population polytopes G, G' , and divergence functionals (e.g., Kullback-Leibler divergence or total variational distance) of the induced marginal data densities. The technical challenges lie in the fact that in order to relate G to the marginal density of the data, one has to integrate out multiple layers of latent variables, $\boldsymbol{\eta}$ and $\boldsymbol{\beta}$. Techniques in convex geometry come in very handily in the derivation of both lower and upper bounds Schneider [1993].

The remainder of the paper is organized as follows. The model and main results are de-

scribed in Section 2. Section 3 describes the basic geometric assumptions and their consequences. An abstract theorem for posterior contraction for $m \times n$ data setting is formulated in Section 4, whose conditions are verified in the subsequent sections. Section 5 proves a contraction result which helps to establish a key lower bound on the Hellinger information, while Section 6 provides a lower bound on Kullback-Leibler neighborhoods of the prior support (that is, a bound the prior thickness). Proofs of main theorems and other technical lemmas are presented in Section 7 and the Appendices.

Notations. $B_p(\boldsymbol{\theta}, r)$ denotes a closed p -dimensional Euclidean ball centered at $\boldsymbol{\theta}$ and has radius r . G_ϵ denotes the Minkowsky sum $G_\epsilon := G + B_{d+1}(\mathbf{0}, \epsilon)$. $\text{bd } G$, $\text{extr } G$, $\text{Diam } G$, $\text{aff } G$, $\text{vol}_p G$ denote the boundary, the set of extreme points, the diameter, the affine span, and the p -dimensional volume of set G , respectively. “Extreme points” and “vertices” are interchangeable throughout this paper. Set-theoretic difference between two sets is defined as $G \triangle G' = (G \setminus G') \cup (G' \setminus G)$. $N(\epsilon, \mathcal{G}, d_{\mathcal{H}})$ denotes the covering number of \mathcal{G} in Hausdorff metric $d_{\mathcal{H}}$. $D(\epsilon, \mathcal{G}, d_{\mathcal{H}})$ is the packing number of \mathcal{G} in Hausdorff metric. Several divergence measures for probability distributions are employed: $K(p, q)$, $h(p, q)$, $V(p, q)$ denote Kullback-Leibler divergence, Hellinger and total variation distance between two densities p and q defined with respect to a measure on a common space: $K(p, q) = \int p \log(p/q)$, $h^2(p, q) = \frac{1}{2} \int (\sqrt{p} - \sqrt{q})^2$ and $V(P, Q) = \frac{1}{2} \int |p - q|$. In addition, we define $K_2 = \int p [\log(p/q)]^2$. Throughout the paper, $f(m, n, \epsilon) \lesssim g(m, n, \epsilon)$, equivalently, $f = O(g)$, means $f(m, n, \epsilon) \leq Cg(m, n, \epsilon)$ for some constant C independent of asymptotic quantities m, n and ϵ – details about the dependence of C are made explicit unless obvious from the context. Similarly, $f(m, n, \epsilon) \gtrsim g(m, n, \epsilon)$ or $f = \Omega(g)$ means $f(m, n, \epsilon) \geq Cg(m, n, \epsilon)$.

2 Main results

Model description. As mentioned in the introduction, the central objects of the admixture model are *population structure* variables $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)$, whose convex hull is called the *population polytope*: $G = \text{conv}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)$. $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k$ reside in d -dimensional probability simplex Δ^d . $k < \infty$ is assumed known. Note that G has at most k vertices (i.e. extreme points) among $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k$.

A random vector $\boldsymbol{\eta} \in G$ is parameterized by $\boldsymbol{\eta} = \beta_1 \boldsymbol{\theta}_1 + \dots, \beta_k \boldsymbol{\theta}_k$, where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k) \in \Delta^{k-1}$ is a random vector distributed according to a distribution $P_{\boldsymbol{\beta}|\gamma}$ for some parameter γ (both Pritchard et al. [2000] and Blei et al. [2003] used the Dirichlet distribution). Given $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k$, this induces a probability distribution $P_{\boldsymbol{\eta}|G}$ whose support is the convex set G . Details of this distribution, suppressed for the time being, are given explicitly by Eq. (14) and (15). [To be precise $P_{\boldsymbol{\eta}|G}$ should be written as $P_{\boldsymbol{\eta}|\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k; G}$. That is, G is always attached with a specific set of $\boldsymbol{\theta}_j$ ’s. Throughout the paper, this specification of G is always understood but notationally suppressed to avoid cluttering.]

For each individual $i = 1, \dots, m$, let $\boldsymbol{\eta}_i \in \Delta^d$ be an independent random vector distributed by $P_{\boldsymbol{\eta}|G}$. The observed data associated with i , $\mathcal{S}_{[n]}^i = (X_{ij})_{j=1}^n$ are assumed to be

i.i.d. draws from the multinomial distribution $\text{Mult}(\boldsymbol{\eta}_i)$ specified by $\boldsymbol{\eta}_i := (\eta_{i0}, \dots, \eta_{id})$. That is, $X_{ij} \in \{0, \dots, d\}$ such that $P(X_{ij} = l | \boldsymbol{\eta}_i) = \eta_{il}$ for $l = 0, \dots, d$.

Admixture models are simple when specified in a hierarchical manner as given above. The relevant distributions are written down below. The joint distribution of the generic random variable $\boldsymbol{\eta}$ and n -vector $\mathcal{S}_{[n]}$ (dropping superscript i used for indexing a specific individual) is denoted by $P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]} | G}$ and its density $p_{\boldsymbol{\eta} \times \mathcal{S}_{[n]} | G}$. We have

$$p_{\boldsymbol{\eta} \times \mathcal{S}_{[n]} | G}(\boldsymbol{\eta}_i, \mathcal{S}_{[n]}^i) = p_{\boldsymbol{\eta} | G}(\boldsymbol{\eta}_i) \times \prod_{j=1}^n \prod_{l=0}^d \eta_{il}^{\mathbb{I}(X_{ij}=l)}. \quad (1)$$

The distribution of $\mathcal{S}_{[n]}$, denoted by $P_{\mathcal{S}_{[n]} | G}$, is obtained by integrating out $\boldsymbol{\eta}$, which yields the following density with respect to counting measure:

$$p_{\mathcal{S}_{[n]} | G}(\mathcal{S}_{[n]}^i) = \int_G \prod_{j=1}^n \prod_{l=0}^d \eta_{il}^{\mathbb{I}(X_{ij}=l)} dP_{\boldsymbol{\eta} | G}(\boldsymbol{\eta}_i). \quad (2)$$

The joint distribution of the full data set $\mathcal{S}_{[n]}^{[m]} := (\mathcal{S}_{[n]}^i)_{i=1}^m$, denoted by $P_{\mathcal{S}_{[n]} | G}^m$, is a product distribution:

$$P_{\mathcal{S}_{[n]} | G}^m(\mathcal{S}_{[n]}^{[m]}) := \prod_{i=1}^m P_{\mathcal{S}_{[n]} | G}(\mathcal{S}_{[n]}^i). \quad (3)$$

Admixture models are customarily introduced in an equivalent way as follows Blei et al. [2003], Pritchard et al. [2000]: For each $i = 1, \dots, m$, draw an independent random variable $\boldsymbol{\beta} \in \Delta^{k-1}$ as $\boldsymbol{\beta} \sim P_{\boldsymbol{\beta} | \gamma}$. Given i and $\boldsymbol{\beta}$, for $j = 1, \dots, n$, draw $Z_{ij} | \boldsymbol{\beta} \stackrel{iid}{\sim} \text{Mult}(\boldsymbol{\beta})$. Z_{ij} takes values in $\{1, \dots, k\}$. Now, data point X_{ij} is randomly generated by $X_{ij} | Z_{ij} = l, \boldsymbol{\theta} \sim \text{Mult}(\boldsymbol{\theta}_l)$. This yields the same joint distribution of $\mathcal{S}_{[n]}^i = (X_{ij})_{j=1}^n$ as the one described earlier. The use of latent variables Z_{ij} is amenable to the development of computational algorithms for inference. However, this representation bears no significance within the scope of this work.

Asymptotic setting and metrics on population polytopes. Assume the data set $\mathcal{S}_{[n]}^{[m]} = (\mathcal{S}_{[n]}^i)_{i=1}^m$ of size $m \times n$ is generated according an admixture model given by “true” parameters $\boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_k^*$. $G_0 = \text{conv}(\boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_k^*)$ is the true population polytope. Under the Bayesian estimation framework, the population structure variables $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)$ are random and endowed with a prior distribution Π . The main question to be addressed in this paper is the contraction behavior of the posterior distribution $\Pi(G | \mathcal{S}_{[n]}^{[m]})$, as the number of data points $m \times n$ goes to infinity.

It is noted that we do not always assume that the number of extreme points of the population polytope G_0 is k . We work in a general overfitted setting where k only serves as the upper bound of the true number of extreme points for the purpose of model parameterization. The special case in which the number of extreme points of G_0 is known a priori is also interesting and will be considered.

Let $\text{extr } G$ denote the set of extreme points of a given polytope G . \mathcal{G}^k is the set of population polytopes in Δ^d such that $|\text{extr } G| \leq k$. Let $\mathcal{G}^* = \cup_{2 \leq k < \infty} \mathcal{G}^k$ be the set of population polytopes that have finite number of extreme points in Δ^d . A natural metric on \mathcal{G}^* is the following “minimum-matching” Euclidean distance:

$$d_{\mathcal{M}}(G, G') = \max_{\theta \in \text{extr } G} \min_{\theta' \in \text{extr } G'} \|\theta - \theta'\| \vee \max_{\theta' \in \text{extr } G'} \min_{\theta \in \text{extr } G} \|\theta' - \theta\|.$$

A more common metric is the Hausdorff metric:

$$d_{\mathcal{H}}(G, G') = \min\{\epsilon \geq 0 \mid G \subset G'_\epsilon; G' \subset G_\epsilon\} = \max_{\theta \in G} d(\theta, G') \vee \max_{\theta' \in G'} d(\theta', G).$$

Here, $G_\epsilon = G + B_{d+1}(\mathbf{0}, \epsilon) := \{\theta + e \mid \theta \in G, e \in \mathbb{R}^{d+1}, \|e\| \leq 1\}$, and $d(\theta, G') := \inf\{\|\theta - \theta'\|, \theta' \in G'\}$. Observe that $d_{\mathcal{H}}$ depends on the boundary structure of sets, while $d_{\mathcal{M}}$ depends on only extreme points. In general, $d_{\mathcal{M}}$ dominates $d_{\mathcal{H}}$, but under additional mild assumptions the two metrics are equivalent (see Lemma 1).

We introduce a notion of regularity for a family probability distributions defined on convex polytopes $G \in \mathcal{G}^*$. This notion is concerned with the behavior near the boundary of the support of distributions $P_{\eta|G}$. We say a family of distributions $\{P_{\eta|G} \mid G \in \mathcal{G}^k\}$ is α -regular if for any $G \in \mathcal{G}^k$ and any $\eta_0 \in \text{bd } G$,

$$P_{\eta|G}(\|\eta - \eta_0\| \leq \epsilon) \geq c\epsilon^\alpha \text{vol}_p(G \cap B_{d+1}(\eta_0, \epsilon)).$$

where p is the number of dimensions of the affine space $\text{aff } G$ that spans G , constant $c > 0$ is independent of G , η_0 and ϵ .

Assumptions. Π is a prior distribution on $\theta_1, \dots, \theta_k$ such that the following hold for the relevant parameters that reside in the support of Π :

- (S0) Geometric properties (A1) and (A2) listed in Section 3 are satisfied uniformly for all G .
- (S1) Each of $\theta_1, \dots, \theta_k$ is bounded away from the boundary of Δ^d . That is, if $\theta_j = (\theta_{j,0}, \dots, \theta_{j,d})$ then $\min_{l=0, \dots, d} \theta_{j,l} > c_0$ for all $j = 1, \dots, k$.
- (S2) For any small ϵ , $\Pi(\|\theta_j - \theta_j^*\| \leq \epsilon \mid \forall j = 1, \dots, k) \geq c'_0 \epsilon^{kd}$, for some $c'_0 > 0$.
- (S3) $\beta = (\beta_1, \dots, \beta_k)$ is distributed (a priori) according to a symmetric probability distribution P_β on Δ^{k-1} . That is, the random variables β_1, \dots, β_k are exchangeable.
- (S4) P_β induces a family of distributions $\{P_{\eta|G} \mid G \in \mathcal{G}^k\}$ that is α -regular.

Theorem 1. Let $G_0 \in \mathcal{G}^k$ and G_0 is in the support of prior Π . Let $p = (k-1) \wedge d$. Under Assumptions (S0–S4) of the admixture model, as $m \rightarrow \infty$ and $n \rightarrow \infty$ such as $\log \log m \leq \log n = o(m)$, for some sufficiently large constant C independent of m and n ,

$$\Pi(d_{\mathcal{M}}(G_0, G) \geq C\delta_{m,n}|\mathcal{S}_{[n]}^{[m]}) \rightarrow 0 \quad (4)$$

in $P_{S_{[n]}|G_0}^m$ -probability. Here,

$$\delta_{m,n} = \left[\frac{\log m}{m} + \frac{\log n}{n} + \frac{\log n}{m} \right]^{\frac{1}{2(p+\alpha)}}.$$

The same statement holds for the Hausdorff metric $d_{\mathcal{H}}$.

Remarks. (i) Geometric assumption (S0) and its consequences are presented in the next section. (S0)(S1) and (S2) are mild assumptions observed in practice (cf. Blei et al. [2003], Pritchard et al. [2000]). (S4) is a standard assumption that holds for a range of α , when $P_{\beta|\gamma}$ is a Dirichlet distribution (see Lemma 4), but there may be other choices. The assumption in (S3) that P_{β} is symmetric is relatively strong, but again it has been widely adopted (e.g., symmetric Dirichlet distributions, including the uniform distribution). It may be difficult to try to relax this assumption if one insists on using Hausdorff metric, see the remark following the statement of Lemma 7.

(ii) In practice P_{β} may be further parameterized as $P_{\beta|\gamma}$, where γ is endowed with a prior distribution. Then, it would be of interest to also study the posterior contraction behavior for γ . In this paper we have opted to focus only on convergence behavior of the population structure to simplify the exposition and the results.

(iii) The appearance of both m^{-1} and n^{-1} in the contraction rate suggests that if either m or n is small, the rate would suffer even if the total amount of data $m \times n$ increases. What is quite interesting is the appearance of $\log n/m$. This is rooted in a condition that the thickness of the prior support in Kullback-Leibler neighborhood for marginal density $p_{S_{[n]}|G}$ is appropriately bounded from below. See Theorem 6 for such a lower bound. This in turn arises from an upper bound on Kullback-Leibler distance of marginal densities, which increases with n (see Lemma 7). From a hierarchical modeling viewpoint, this result highlights an interesting interaction of sample sizes provided to different levels in the model hierarchy. This issue has not been widely discussed in the hierarchical modeling literature in a theoretical manner, to the best of our knowledge.

(iv) Note the constraints that $n > \log m$ and $\log n = o(m)$ are required in order to obtain rates of posterior contraction. These constraints are related to the term $\log n/m$ mentioned above — they stem from the upper bound on Kullback-Leibler in Lemma 7. The remark following the statement of this lemma explains why the upper bound almost always grows with n . A very special situation is presented in Lemma 5 where an upper bound on Kullback-Leibler distance can be obtained that is independent of n . However, such a situation cannot be verified in any reasonable estimation setting. This suggests that with our proof technique, we almost always require n to grow at a constrained rate relatively to m in order to obtain posterior contraction rates.

(v) Since the quantity $\log n/m$ arises partly from a fairly general proof technique in Bayesian asymptotics, one may wonder whether it is possible to get rid of it by considering a point estimation procedure for G . A line of reasoning goes like this. For each $i = 1, \dots, m$, as $n \rightarrow \infty$ one can estimate η_i arbitrarily well at a rate $O(n^{-1/2})$. All that

remain is to estimate polytope G as if the (exact) samples from $P_{\eta|G}$ are available, an estimation task that incurs a rate depending on m , independent of n . Since one does not have exact samples from $P_{\eta|G}$, a more formal reasoning of similar spirit is needed: Observe that all information we have access to regarding polytope G is through marginal density $p_{\mathcal{S}_{[n]}|G}$ for m samples of n -vector $\mathcal{S}_{[n]}$. By Theorem 5: if $h(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G_0}) \rightarrow 0$ at a rate, say ϵ_m , as $m \rightarrow \infty$ and $n \rightarrow \infty$ suitably, then $d_{\mathcal{H}}(G, G_0)$ vanishes at the rate $(\epsilon_m + (\log n/n)^{1/2})^{1/\gamma}$ for some constant $\gamma > 0$. If we can somehow show that a point estimate for $p_{\mathcal{S}_{[n]}|G_0}$ (e.g., the maximum likelihood estimator) can yield a parametric rate in Hellinger metric, say $\epsilon_m \asymp (\log m/m)^{1/2}$, then the convergence rate for G in $d_{\mathcal{H}}$ would be $(\log m/m + \log n/n)^{1/2\gamma}$, which is happily free of $\log n/m$. It is not so obvious if this is possible, as sample size n should nonetheless affect the complexity of random vectors $\mathcal{S}_{[n]}$. Formally, the entropy number of the space of marginal densities $p_{\mathcal{S}_{[n]}|G}$ generally scales with $O(\log n)$. A standard derivation (see, e.g., van de Geer [2000]) would still yield the rate $\epsilon_m \asymp (\log n/m)^{1/2}$. Our conclusion is that without strong assumptions such as the one discussed in remark (iv), removing quantities such as $\log n/m$ from the convergence rate is far from trivial.

(vi) The exponent $\frac{1}{2(p+\alpha)}$ suggests a slow, nonparametric-like convergence rate. Moreover, later in Theorem 3 we show that this is qualitatively quite close to a minimax lower bound. On the other hand, the following theorem shows that it is possible to achieve a parametric rate if additional constraints are imposed on the true G_0 and/or the prior Π :

Theorem 2. *Let $G_0 \in \mathcal{G}^k$ and G_0 is in the support of prior Π . Assume (S0–S4), and either one of the following two conditions hold:*

- (a) $|\text{extr } G_0| = k$, or
- (b) *There is a known constant $r_0 > 0$ such that the pairwise distances of the extreme points of all G in the support of the prior are bounded from below by r_0 .*

Then, as $m \rightarrow \infty$ and $n \rightarrow \infty$ such that $\log m < n$ and $\log n = o(m)$, Eq. (4) holds with

$$\delta_{m,n} = \left[\frac{\log m}{m} + \frac{\log n}{n} + \frac{\log n}{m} \right]^{\frac{1}{2(1+\alpha)}},$$

The same statement holds for the Hausdorff metric $d_{\mathcal{H}}$.

The next theorem produces minimax lower bounds that are qualitatively quite close to the nonparametric-like rates obtained in Theorem 1. In the following theorem, η is not parameterized by β and θ_j 's as in the admixture model. Instead, we shall simply replace assumptions (S3) and (S4) on $P_{\beta|\gamma}$ by either one of the following assumptions on $P_{\eta|G}$:

- (S5) For any pair of p -dimensional polytopes $G' \subset G$ that satisfy property A1,

$$V(P_{\eta|G}, P_{\eta|G'}) \lesssim d_{\mathcal{H}}(G, G')^\alpha \text{vol}_p G \setminus G'.$$

(S5') For any p -dimensional polytope G , $P_{\eta|G}$ is the uniform distribution on G . (This actually entails (S5) for $\alpha = 0$.)

Since a parameterization for η is not needed, the overall model can be simplified as follows: Given population polytope $G \in \Delta^d$, for each $i = 1, \dots, m$, draw $\eta_i \stackrel{iid}{\sim} P_{\eta|G}$. For each $j = 1, \dots, n$, draw $\mathcal{S}_{[n]}^i = (X_{ij})_{j=1}^n \stackrel{iid}{\sim} \text{Mult}(\eta_i)$.

Theorem 3. Suppose that $G_0 \in \mathcal{G}^k$ satisfies assumptions (S0)(S1) and (S2). Point estimates $\hat{G} = \hat{G}(\mathcal{S}_{[n]}^{[m]})$ are restricted to those that also satisfy (S0)(S1) and (S2). In the following, infimum and supremum are taken over the specified domains for G_0 and \hat{G} , while the multiplying constants in \gtrsim depend only on constants specified by these assumptions.

(a) Let $q = \lfloor k/2 \rfloor \wedge d$. Under Assumption (S5), we have

$$\inf_{\hat{G}} \sup_{G_0} P_{\mathcal{S}_{[n]}|G_0}^m d_{\mathcal{H}}(G_0, \hat{G}) \gtrsim \left(\frac{1}{mn} \right)^{\frac{1}{q+\alpha}}.$$

(b) Let $q = \lfloor k/2 \rfloor \wedge d$. Under Assumption (S5'), we have

$$\inf_{\hat{G}} \sup_{G_0} P_{\mathcal{S}_{[n]}|G_0}^m d_{\mathcal{H}}(G_0, \hat{G}) \gtrsim \left(\frac{1}{m} \right)^{\frac{1}{q}}.$$

(c) Assume (S5'), and that either condition (a) or (b) of Theorem 2 holds, then

$$\inf_{\hat{G}} \sup_{G_0} P_{\mathcal{S}_{[n]}|G_0}^m d_{\mathcal{H}}(G_0, \hat{G}) \gtrsim \left(\frac{1}{mn} \right)^{\frac{1}{1+\alpha}}.$$

Furthermore, if (S5) is replaced by (S5'), the lower bound becomes $1/m$.

Remarks. (i) Although there remain some gap between the posterior contraction rate in Theorem 1 and the minimax lower bound in Theorem 3 (a), they are qualitatively comparable and both notably dependent on d and k . The gap should be expected partly because slightly enlarged models are considered in Theorem 3, due to the relaxed parameterization. Nonetheless, if $k \geq 2d$, and allowing $m \asymp n$, the rate exponents differ by only a factor of 4. That is, $m^{-1/2(d+\alpha)}$ vis-à-vis $m^{-2/(d+\alpha)}$.

(ii) The nonparametrics-like lower bounds in part (a) and (b) in the overfitted setting are somewhat surprising even if P_{β} is known exactly (e.g., P_{β} is uniform distribution). Since we are more likely to be in the overfitted setting than knowing the exact number of extreme points, an implication of this is that it is important to in practice to impose a lower bound on the pairwise distances between the extreme points of the population polytope.

(iii) The results in part (b) and (c) under assumption (S5') present an interesting scenario in which the obtained lower bounds do not depend on n , which determines the amount of data at the bottom level in the model hierarchy.

3 Geometric assumptions and basic lemmas

In this section we discuss the geometric assumptions postulated in the main theorems, and describe their consequences using elementary arguments in convex geometry of Euclidean spaces. These results relate Hausdorff metric, the minimum-matching metric, and the volume of the set-theoretic difference of polytopes. These relationships prove crucial in obtaining explicit posterior contraction rates. Here, we state the properties and prove the results for p -dimensional polytopes and convex bodies of points in Δ^d , for a given $p \leq d$. (Convex bodies are bounded convex sets that may have an unbounded number of extreme points. Within this section, the detail of the ambient space is irrelevant. For instance, Δ^d may be replaced by \mathbb{R}^{d+1} or a higher dimensional Euclidean space).

Property A1. (Property of thick body): For some $r, R > 0$, $\theta_c \in \Delta^d$, G contains the spherical ball $B_p(\theta_c, r)$ and is contained in $B_p(\theta_c, R)$.

Property A2. (Property of non-obstute corners): For some small $\delta > 0$, at each vertex of G there is a supporting hyperplane whose angle formed with any edges adjacent to that vertex is bounded from below by δ .

We state key geometric lemmas that will be used throughout the paper. Bounds such as those given by Lemma 2 are probably well-known in the folklore of convex geometry (for instance, part (b) of that lemma is similar to (but not precisely the same as) Lemma 2.3.6. from Schneider [1993]). Due to the absence of direct references we include the proof of this and other lemmas in the Appendix.

Lemma 1. (a) $d_{\mathcal{H}}(G, G') \leq d_{\mathcal{M}}(G, G')$.

(b) If the two polytopes G, G' satisfy property A2, then $d_{\mathcal{M}}(G, G') \leq C_0 d_{\mathcal{H}}(G, G')$, for some positive constant $C_0 > 0$ depending only on δ .

According to part (b) of this lemma, convergence of a sequence of convex polytope $G \in \mathcal{G}^k$ to $G_0 \in \mathcal{G}^k$ in Hausdorff metric entails the convergence of the extreme points of G to those of G_0 . Moreover, they share the same rate as the Hausdorff convergence.

Lemma 2. There are positive constants C_1 and c_1 depending only on r, R, p such that for any two p -dimensional convex bodies G, G' satisfying property A1:

$$(a) \text{ vol}_p G \triangle G' \geq c_1 d_{\mathcal{H}}(G, G')^p.$$

$$(b) \text{ vol}_p G \triangle G' \leq C_1 d_{\mathcal{H}}(G, G').$$

Remark. The exponents in both bounds in Lemma 2 are attainable. Indeed, for the lower bound in part (a), consider a fixed convex polytope G . For each vertex $\theta_i \in G$, consider point x that lie on edges incident to θ_i such that $\|x - \theta_i\| = \epsilon$. Let G' be the convex hull of all such x 's and the remaining vertices of G . Clearly, $d_{\mathcal{H}}(G, G') = O(\epsilon)$, and $\text{vol}_p G \setminus G' \leq O(\epsilon^p)$. Thus, for the collection of convex polytopes G' constructed in this

way, $\text{vol}_p(G \triangle G') \asymp d_{\mathcal{H}}(G, G')^p$. The upper bound in part (b) is also tight for a broad class of convex polytopes, as exemplified by the following lemma.

Lemma 3. *Let G be a fixed polytope and $|\text{extr } G| = k < \infty$. G' an arbitrary polytope in \mathcal{G}^* . Moreover, either one of the following conditions holds:*

- (a) $|\text{extr } G'| = k$, or
- (b) *The pairwise distances between the extreme points of G' is bounded away from a constant $r_0 > 0$.*

Then, there is a positive constant $\epsilon_0 = \epsilon_0(G)$ depending only on G , a positive constant $c_2 = c_2(G)$ in case (a) and $c_2 = c_2(G, r_0)$ in case (b), such that

$$\text{vol}_p G \triangle G' \geq c_2 d_{\mathcal{H}}(G, G')$$

as soon as $d_{\mathcal{H}}(G, G') \leq \epsilon_0(G)$.

Remark. We note that the bound obtained in this lemma is substantially stronger than the lemma obtained by Lemma 2 part (a). This is due to the asymmetric roles of G , which is held fixed, and G' , which can vary. As a result, constant c_2 as stated in the present lemma is independent of G' but allowed to be dependent on G . By contrast, constant c_1 in Lemma 2 part (a) is independent of both G and G' .

4 An abstract posterior contraction theorem

In this section we state an abstract posterior contraction theorem for hierarchical models, whose proof is given in the Appendix. The setting of this theorem is a general hierarchical model defined as follows

$$G \sim \Pi, \boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_m | G \sim P_{\boldsymbol{\eta}|G}$$

$$\mathcal{S}_{[n]}^i | \boldsymbol{\eta}_i \sim P_{\mathcal{S}_{[n]} | \boldsymbol{\eta}_i} \text{ for } i = 1, \dots, m.$$

The detail of conditional distributions in above specifications is actually irrelevant. Thus results in this section may be of general interest for hierarchical models with $m \times n$ data.

As before $p_{\mathcal{S}_{[n]}|G}$ is marginal density of the generic $\mathcal{S}_{[n]}$ which is obtained by integrating out the generic random vector $\boldsymbol{\eta}$ (e.g., see Eq. (2)). We need several key notions. Define the Hausdorff ball as:

$$B_{d_{\mathcal{H}}}(G_1, \delta) := \{G \in \Delta^d : d_{\mathcal{H}}(G_1, G) \leq \delta\}.$$

A useful quantity for proving posterior concentration theorems is the Hellinger information of Hausdorff metric for a given set:

Definition 1. Fix $G_0 \in \mathcal{G}^k$. \mathcal{G} is a subset of \mathcal{G}^k . For a fixed n , the sample size of $\mathcal{S}_{[n]}$, define the Hellinger information of $d_{\mathcal{H}}$ metric for set \mathcal{G} as a real-valued function on the positive reals $\Psi_{\mathcal{G},n} : \mathbb{R}_+ \rightarrow \mathbb{R}$:

$$\Psi_{\mathcal{G},n}(\delta) := \inf_{G \in \mathcal{G}; d_{\mathcal{H}}(G_0, G) \geq \delta/2} h^2(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G}). \quad (5)$$

We also define $\Phi_{\mathcal{G},n} : \mathbb{R}_+ \rightarrow \mathbb{R}$ to be an arbitrary non-negative valued function on the positive reals such that for any $\delta > 0$,

$$\sup_{G, G' \in \mathcal{G}; d_{\mathcal{H}}(G, G') \leq \Phi_{\mathcal{G},n}(\delta)} h^2(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq \Psi_{\mathcal{G},n}(\delta)/4.$$

In both definitions of Φ and Ψ , we suppress the dependence on (the fixed) G_0 and δ to simplify notations. Note that if $G_0 \in \mathcal{G}$, it follows from the definition that $\Phi_{\mathcal{G},n}(\delta) < \delta/2$.

Remark. Suppose that conditions of Lemma 7 (b) hold, so that

$$h^2(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq \frac{n}{c_0} C_0 d_{\mathcal{H}}(G, G').$$

Then it suffices to choose $\Phi_{\mathcal{G},n}(\delta) = \frac{c_0}{4nC_0} \Psi_{\mathcal{G},n}(\delta)$.

Define the neighborhood of the prior support around G_0 in terms of Kullback-Leibler distance of the marginal densities $p_{\mathcal{S}_{[n]}|G}$:

$$B_K(G_0, \delta) = \{G \in \mathcal{G}^* | K(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G}) \leq \delta^2; K_2(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G}) \leq \delta^2\}. \quad (6)$$

Theorem 4. Suppose that

- (a) $m \rightarrow \infty$ and $n \rightarrow \infty$ at a certain rate relative to m ,
- (b) There is a sequence of subsets $\mathcal{G}_m \subset \mathcal{G}^*$, a large constant C , a sequence of scalars $\epsilon_{m,n} \rightarrow 0$ defined in terms of m and n such that $m\epsilon_{m,n}^2$ tends to infinity, such that

$$\sup_{G_1 \in \mathcal{G}_m} \log D(\Phi_{\mathcal{G}_m,n}(\epsilon), \mathcal{G}_m \cap B_{d_{\mathcal{H}}}(G_1, \epsilon/2), d_{\mathcal{H}}) \quad (7)$$

$$+ \log D(\epsilon/2, \mathcal{G}_m \cap B_{d_{\mathcal{H}}}(G_0, 2\epsilon) \setminus B_{d_{\mathcal{H}}}(G_0, \epsilon), d_{\mathcal{H}}) \leq m\epsilon_{m,n}^2$$

$$\forall \epsilon \geq \epsilon_{m,n},$$

$$\Pi(\mathcal{G}^* \setminus \mathcal{G}_m) \leq \exp[-m\epsilon_{m,n}^2(C+4)], \quad (8)$$

$$\Pi(B_K(G_0, \epsilon_{m,n})) \geq \exp[-m\epsilon_{m,n}^2 C]. \quad (9)$$

- (c) There is a sequence of positive scalars M_m such that

$$\Psi_{\mathcal{G}_m,n}(M_m \epsilon_{m,n}) \geq 8\epsilon_{m,n}^2(C+4) \quad (10)$$

$$\exp(2m\epsilon_{m,n}^2) \sum_{j \geq M_m} \exp[-m\Psi_{\mathcal{G}_m,n}(j\epsilon_{m,n})/8] \rightarrow 0. \quad (11)$$

Then, $\Pi(G : d_{\mathcal{H}}(G_0, G) \geq M_m \epsilon_{m,n} |\mathcal{S}_{[n]}^{[m]}) \rightarrow 0$ in $P_{\mathcal{S}_{[n]}|G_0}^m$ -probability as m and $n \rightarrow \infty$.

The proof of this theorem is deferred to the Appendix. As noted above, this result is applicable to any hierarchical models for $m \times n$ data. The choice of Hausdorff metric $d_{\mathcal{H}}$ is arbitrary here, and can be replaced by any other valid metric (e.g., $d_{\mathcal{M}}$). The remainder of the paper is devoted to verifying the conditions of this theorem so it can be applied. These conditions hinge on our having established a lower bound for the Hellinger information function $\Psi_{\mathcal{G}_m, n}(\cdot)$ (via Theorem 5), and a lower bound for the prior probability defined on Kullback-Leibler balls $B_K(G_0, \cdot)$ (via Theorem 6). Both types of results are obtained by utilizing the convex geometry lemmas described in the previous section.

5 Contraction properties

The following contraction result guarantees that as marginal densities of $\mathcal{S}_{[n]}$ get closer in total variation distance metric (or Hellinger metric), so do the corresponding population polytopes in Hausdorff metric (or minimum matching metric). This gives a lower bound for the Hellinger information defined by Eq. (5), because h is related to V via inequality $h \geq V$.

Theorem 5. (a) Let G, G' be two convex bodies in Δ^d . G is a p -dimensional body containing spherical ball $B_p(\theta_c, r)$, while G' is p' -dimensional body containing $B_{p'}(\theta_c, r)$ for some $p, p' \leq d, r > 0, \theta_c \in \Delta^d$. In addition, assume that both $p_{\eta|G}$ and $p_{\eta|G'}$ are α -regular densities on G and G' , respectively. Then, there is $c_1 > 0$ independent of G, G' such that

$$c_1 d_{\mathcal{H}}(G, G')^{(p \vee p') + \alpha} \leq V(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) + 6(d+1) \exp \left[- \frac{n}{8(d+1)} d_{\mathcal{H}}(G, G')^2 \right].$$

(b) Assume further that G is fixed convex polytope, G' an arbitrary polytope, $p' = p$, and that either $|\text{extr } G'| = |\text{extr } G|$ or the pairwise distances of extreme points of G' is bounded from below by a constant $r_0 > 0$. Then, there are constants $c_2, C_3 > 0$ depending only on G and r_0 (and independent of G') such that

$$c_2 d_{\mathcal{H}}(G, G')^{1+\alpha} \leq V(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) + 6(d+1) \exp \left[- \frac{n}{C_3(d+1)} d_{\mathcal{H}}(G, G')^2 \right].$$

Remark. Part (a) holds for varying pairs of G, G' satisfying certain conditions. It is consequence of Lemma 2 (a). Part (b) produces a tighter bound, but it holds only for a fixed G , while G' is allowed to vary while satisfying certain conditions. This is a consequence of Lemma 3. Constants c_1, c_2 are the same as those from Lemma 2 (a) and 3, respectively.

Proof. (a) The main idea of the proof is the construction of a suitable test set in order to distinguish $p_{\mathcal{S}_{[n]}|G'}$ from $p_{\mathcal{S}_{[n]}|G}$. The proof is organized as a sequence of steps.

Step 1 Given a data vector $\mathcal{S}_{[n]} = (X_1, \dots, X_n)$, define $\hat{\boldsymbol{\eta}}(\mathcal{S}) \in \Delta^d$ such that the i -element of $\hat{\boldsymbol{\eta}}(\mathcal{S})$ is $\frac{1}{n} \sum_{j=1}^n \mathbb{I}(X_j = i)$ for each $i = 0, \dots, d$. In the following we simply use $\hat{\boldsymbol{\eta}}$ to ease the notations. By the definition of the variational distance,

$$V(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) = \sup_A |P_{\mathcal{S}_{[n]}|G}(\hat{\boldsymbol{\eta}} \in A) - P_{\mathcal{S}_{[n]}|G'}(\hat{\boldsymbol{\eta}} \in A)|, \quad (12)$$

where the supremum is taken over all measurable subsets of Δ^d .

Step 2 Fix a constant $\epsilon > 0$. By Hoeffding's inequality and the union bound, under the conditional distribution $P_{\mathcal{S}_{[n]}|\boldsymbol{\eta}}$,

$$P_{\mathcal{S}_{[n]}|\boldsymbol{\eta}}(\max_{i=0,\dots,d} |\hat{\eta}_i - \eta_i| \geq \epsilon) \leq 2(d+1) \exp(-2n\epsilon^2)$$

with probability one (as $\boldsymbol{\eta}$ is random). It follows that

$$\begin{aligned} P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G}(\|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}\| \geq \epsilon) &\leq P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G}(\max_{i=0,\dots,d} |\hat{\eta}_i - \eta_i| \geq \epsilon(d+1)^{-1/2}) \\ &\leq 2(d+1) \exp[-2n\epsilon^2/(d+1)]. \end{aligned}$$

The same bound holds under $P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G'}$.

Step 3 Define event $B = \{\|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}\| < \epsilon\}$. Take any (measurable) set $A \subset \Delta^d$,

$$\begin{aligned} &|P_{\mathcal{S}_{[n]}|G}(\hat{\boldsymbol{\eta}} \in A) - P_{\mathcal{S}_{[n]}|G'}(\hat{\boldsymbol{\eta}} \in A)| \\ &= |P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G}(\hat{\boldsymbol{\eta}} \in A; B) + P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G}(\hat{\boldsymbol{\eta}} \in A; B^C) \\ &\quad - P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G'}(\hat{\boldsymbol{\eta}} \in A; B) - P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G'}(\hat{\boldsymbol{\eta}} \in A; B^C)| \\ &\geq |P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G}(\hat{\boldsymbol{\eta}} \in A; B) - P_{\boldsymbol{\eta} \times \mathcal{S}_{[n]}|G'}(\hat{\boldsymbol{\eta}} \in A; B)| \\ &\quad - 4(d+1) \exp[-2n\epsilon^2/(d+1)]. \end{aligned} \quad (13)$$

Step 4 Let $\epsilon_1 = d_{\mathcal{H}}(G, G')/4$. For any $\epsilon \leq \epsilon_1$, recall the outer ϵ -parallel set $G_\epsilon = (G + B_{d+1}(\mathbf{0}, \epsilon))$, which is full-dimensional $(d+1)$ eventhough G may not be. By triangular inequality, $d_{\mathcal{H}}(G_\epsilon, G'_\epsilon) \geq d_{\mathcal{H}}(G, G')/2$. We shall argue that for any $\epsilon \leq \epsilon_1$, there is a constant $c_1 > 0$ independent of G, G', ϵ and ϵ_1 such that either one of the two scenarios holds:

- (i) There is a set $A^* \subset G \setminus G'$ such that $A_\epsilon^* \cap G'_\epsilon = \emptyset$ and $\text{vol}_p(A^*) \geq c_1 \epsilon_1^p$, or
- (ii) There is a set $A^* \subset G' \setminus G$ such that $A_\epsilon^* \cap G_\epsilon = \emptyset$ and $\text{vol}_{p'}(A^*) \geq c_1 \epsilon_1^{p'}$.

Indeed, since $\epsilon \leq d_{\mathcal{H}}(G, G')/4$, either one of the following two inequalities holds: $d_{\mathcal{H}}(G \setminus G'_{3\epsilon}, G') \geq d_{\mathcal{H}}(G, G')/4$ or $d_{\mathcal{H}}(G' \setminus G_{3\epsilon}, G) \geq d_{\mathcal{H}}(G, G')/4$. If the former inequality holds, let $A^* = G \setminus G'_{3\epsilon}$. Then, $A^* \subset G \setminus G'$ and $A_\epsilon^* \cap G'_\epsilon = \emptyset$. Moreover, by

Lemma 2 (a), $\text{vol}_p(A^*) \geq c_1 \epsilon_1^p$, for some constant $c_1 > 0$ independent of $\epsilon, \epsilon_1, G, G'$, so A^* satisfies (i). In fact, using the same argument as in the proof of Lemma 2 (a) there is a point $x \in \text{bd } G$ such that $G' \cap B_p(x, \epsilon_1) = \emptyset$. Combined with the α -regularity of $P_{\eta|G}$, we have $P_{\eta|G}(A^*) \geq \epsilon^\alpha \text{vol}_p(G \cap B_p(x, \epsilon_1)) \geq c_1 \epsilon^{p+\alpha}$ for some constant $c_1 > 0$. If the latter inequality holds, the same argument applies by defining $A^* = G' \setminus G_{3\epsilon}$ so that (ii) holds.

Step 5 Suppose that (i) holds for the chosen A^* . This means that $P_{\eta \times \mathcal{S}_{[n]}|G'}(\hat{\eta} \in A_\epsilon^*; B) \leq P_{\eta|G'}(\eta \in A_{2\epsilon}^*) = 0$, since $A_{2\epsilon}^* \cap G' = \emptyset$, which is a consequence of $A_\epsilon^* \cap G'_\epsilon = \emptyset$. In addition,

$$\begin{aligned} P_{\eta \times \mathcal{S}_{[n]}|G}(\hat{\eta} \in A_\epsilon^*; B) &\geq P_{\eta \times \mathcal{S}_{[n]}|G}(\eta \in A^*; B) \\ &\geq P_{\eta|G}(A^*) - P_{\eta \times \mathcal{S}_{[n]}|G}(B^C) \\ &\geq P_{\eta|G}(A^*) - 2(d+1) \exp(-2n\epsilon^2/(d+1)) \\ &\geq c_1 \epsilon_1^{p+\alpha} - 2(d+1) \exp(-2n\epsilon^2/(d+1)). \end{aligned}$$

Hence, by Eq. (13) $|P_{\mathcal{S}_{[n]}|G}(\hat{\eta} \in A_\epsilon^*) - P_{\mathcal{S}_{[n]}|G'}(\hat{\eta} \in A_\epsilon^*)| \geq c_1 \epsilon_1^{p+\alpha} - 6(d+1) \exp(-2n\epsilon^2)$. Set $\epsilon = \epsilon_1$, the conclusion then follows by invoking Eq. (12). The scenario of (ii) proceeds in the same way.

(b) Under the condition that the pairwise distances of extreme points of G' are bounded from below by $r_0 > 0$, the proof is very similar to part (a), by involving Lemma 3. Under the condition that $|\text{extr } G'| = k$, the proof is also similar, but it requires a suitable modification for the existence of set A^* . For any small ϵ , let \tilde{G}_ϵ be the minimum-volume homothetic transformation of G , with respect to center θ_c , such that \tilde{G}_ϵ contains G_ϵ . Since $B_p(\theta_c, r) \subset G \subset B_p(\theta_c, R)$ for $R = 1$, it is simple to see that $d_{\mathcal{H}}(G, \tilde{G}_\epsilon) \leq \epsilon R/r = \epsilon/r$.

Set $\epsilon_1 = d_{\mathcal{H}}(G, G')r/4$. We shall argue that for any $\epsilon \leq \epsilon_1$, there is a constant $c_0 > 0$ independent of G', ϵ and ϵ_1 such that either one of the following two scenarios hold:

(iii) There is a set $A^* \subset G \setminus G'$ such that $A_\epsilon^* \cap G'_\epsilon = \emptyset$ and $\text{vol}_p(A^*) \geq c_2 \epsilon_1$, or

(iv) There is a set $A^* \subset G' \setminus G$ such that $A_\epsilon^* \cap G_\epsilon = \emptyset$ and $\text{vol}_p(A^*) \geq c_2 \epsilon_1$.

Indeed, note that either one of the following two inequalities holds: $d_{\mathcal{H}}(G \setminus \tilde{G}'_{3\epsilon}, G') \geq d_{\mathcal{H}}(G, G')/4$ or $d_{\mathcal{H}}(G' \setminus \tilde{G}_{3\epsilon}, G) \geq d_{\mathcal{H}}(G, G')/4$. If the former inequality holds, let $A^* = G \setminus \tilde{G}'_{3\epsilon}$. Then, $A^* \subset G \setminus G'$ and $A_\epsilon^* \cap \tilde{G}'_\epsilon = \emptyset$. Observe that both G and $\tilde{G}'_{3\epsilon}$ have the same number of extreme points by the construction. Moreover, G is fixed so that all geometric properties A2, A1 are satisfied for both G and $\tilde{G}'_{3\epsilon}$ for sufficiently small $d_{\mathcal{H}}(G, G')$. By Lemma 3, $\text{vol}_p(A^*) \geq c_2 \epsilon_1$. Hence, (iii) holds. If the latter inequality holds, the same argument applies by defining $A^* = G' \setminus \tilde{G}_{3\epsilon}$ so that (iv) holds.

Now the proof of the theorem proceeds in the same manner as in part (a). □

6 Concentration properties of the prior support

In this section we study properties of the support of the prior probabilities as specified by the admixture model, including bounds for the support of the prior as defined by Kullback-Leibler neighborhoods.

α -regularity. Let β be a random variable taking values in Δ^{k-1} that has a density p_β (with respect to the $k-1$ -dimensional Hausdorff measure on \mathbb{R}^k). Define random variable $\eta = \beta_1 \theta_1 + \dots + \beta_k \theta_k$, which takes values in $G = \text{conv}(\theta_1, \dots, \theta_k)$. Write $\eta = L\beta$, where $L = [\theta_1 \dots \theta_k]$ is a $(d+1) \times k$ matrix. If $k \leq d+1$, $\theta_1, \dots, \theta_k$ are generally linearly independent, in which case matrix L has rank $k-1$. By the change of variable formula Evans and Gariepy [1992] (Chapter 3), P_β induces a distribution $P_{\eta|G}$ on $G \subset \Delta^d$, which admits the following density with respect to the $k-1$ dimensional Hausdorff measure on Δ^d :

$$p_\eta(\eta|G) = p_\beta(L^{-1}(\eta))J(L)^{-1}. \quad (14)$$

Here $J(L)$ denotes the Jacobian of the linear map. On the other hand, if $k \geq d+1$, then L is generally d -ranked. The induced distribution for η admits the following density with respect to the d -dimensional Hausdorff measure on \mathbb{R}^{d+1} :

$$p_\eta(\eta|G) = \int_{L^{-1}\{\eta\}} p_\beta(\beta) J(L)^{-1} \mathcal{H}^{k-(d+1)}(d\beta). \quad (15)$$

A common choice for P_β is the Dirichlet distribution, as adopted by Pritchard et al. [2000], Blei et al. [2003]: given parameter $\gamma \in \mathbb{R}_+^k$, for any $A \subset \Delta^{k-1}$,

$$P_\beta(\beta \in A|\gamma) = \int_A \frac{\Gamma(\sum \gamma_j)}{\prod_{j=1}^k \Gamma(\gamma_j)} \prod_{j=1}^k \beta_j^{\gamma_j-1} \mathcal{H}^{k-1}(d\beta).$$

Lemma 4. Let $\eta = \sum_{j=1}^k \beta_j \theta_j$, where β is distributed according to a $k-1$ -dimensional Dirichlet distribution with parameters $\gamma_j \in (0, 1]$ for $j = 1, \dots, k$.

- (a) If $k \leq d+1$, there is constant $\epsilon_0 = \epsilon_0(k) > 0$, and constant $c_6 = c_6(\gamma, k, d) > 0$ dependent on γ, k and d such that for any $\epsilon < \epsilon_0$,

$$\inf_{G \subset \Delta^d} \inf_{\eta^* \in G} P_{\eta|G}(\|\eta - \eta^*\| \leq \epsilon) \geq c_6 \epsilon^{k-1}.$$

- (b) If $k > d+1$, the statement holds with a lower bound $c_6 \epsilon^{d+\sum_{i=1}^k \gamma_i}$.

A consequence of this lemma is that if $\gamma_j \leq 1$ for all $j = 1, \dots, k$, $k \leq d+1$ and G is $k-1$ -dimensional, then the induced $P_{\eta|G}$ has a Hausdorff density that is bounded away from 0 on the entire its support Δ^{k-1} , which implies 0-regularity. On the other hand, if $\gamma_j \leq 1$ for all j , $k > d+1$, and G is d -dimensional, the $P_{\eta|G}$ is at least $\sum_{j=1}^k \gamma_j$ -regularity. Note that the α -regularity condition is concerned with the density behavior near the boundary of its support, and thus is weaker than what is guaranteed here.

Bounds on KL divergences. Suppose that the population polytope G is endowed with a prior distribution on \mathcal{G}^k (via prior on the population structures $\theta_1, \dots, \theta_k$). Given G , the marginal density $p_{\mathcal{S}_{[n]}|G}$ of n -vector $\mathcal{S}_{[n]}$ is obtained via Eq. (2). To establish the concentration properties of Kullback-Leibler neighborhood B_K as induced by the prior, we need to obtain an upper bound on the KL divergences for the marginal densities in terms of Hausdorff metric on population polytopes. First, consider a very special case:

Lemma 5. *Let $G, G' \in \Delta^d$ be closed convex sets satisfying property A1. Moreover, assume that*

- (a) $G \subset G'$, $\text{aff } G = \text{aff } G'$ is p -dimensional, for $p \leq d$.
- (b) $P_{\eta|G}$ (resp. $P_{\eta|G'}$) are uniform distributions on G , (resp. G').

Then, there is a constant $C_1 = C_1(r, p) > 0$ such that $K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq C_1 d_{\mathcal{H}}(G, G')$.

Proof. First, we note a well-known fact of KL divergences: the divergence between marginal distributions (e.g., on $\mathcal{S}_{[n]}$) is bounded from above by the divergence between joint distributions (e.g., on η and $\mathcal{S}_{[n]}$) via Eq. (1):

$$K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq K(P_{\eta \times \mathcal{S}_{[n]}|G}, P_{\eta \times \mathcal{S}_{[n]}|G'}).$$

Due to the hierarchical specification, $p_{\eta \times \mathcal{S}|G} = p_{\eta|G} \times p_{\mathcal{S}_{[n]}|\eta}$ and $p_{\eta \times \mathcal{S}|G'} = p_{\eta|G'} \times p_{\mathcal{S}_{[n]}|\eta}$, so $K(P_{\eta \times \mathcal{S}_{[n]}|G}, P_{\eta \times \mathcal{S}_{[n]}|G'}) = K(p_{\eta|G}, p_{\eta|G'})$. The assumption $\text{aff } G = \text{aff } G'$ and moreover $G \subset G'$ implies that $K(p_{\eta|G}, p_{\eta|G'}) < \infty$. In addition, $P_{\eta|G}$ and $P_{\eta|G'}$ are assumed to be uniform distributions on G and G' , respectively, so

$$K(p_{\eta|G}, p_{\eta|G'}) = \int \log \frac{1/\text{vol}_p G}{1/\text{vol}_p G'} dP_{\eta|G}.$$

By Lemma 2 (b), $\log[\text{vol}_p G' / \text{vol}_p G] \leq \log(1 + C_1 d_{\mathcal{H}}(G, G')) \leq C_1 d_{\mathcal{H}}(G, G')$ for some constant $C_1 = C_1(r, p) > 0$. This completes the proof. \square

Remark. The previous lemma requires a particular stringent condition, $\text{aff } G = \text{aff } G'$, and moreover $G \subset G'$, which is usually violated when $k < d + 1$. However, the conclusion is worth noting in that the upper bound does not depend on the sample size n (for $\mathcal{S}_{[n]}$). The next lemma removes this condition and the condition that both $p_{\eta|G}$ and $p_{\eta|G'}$ be uniform. As a result the upper bound obtained is weaker, in the sense that the bound is not in terms of a Hausdorff distance, but in terms of a Wasserstein distance.

Let $Q(\eta_1, \eta_2)$ denote a coupling of $P(\eta|G)$ and $P(\eta|G')$, i.e., a joint distribution on $G \times G'$ whose induced marginal distributions of η_1 and η_2 are equal to $P(\eta|G)$ and $P(\eta|G')$, respectively. Let \mathcal{Q} be the set of all such couplings. The Wasserstein distance between $p_{\eta|G}$ and $p_{\eta|G'}$ is defined as

$$W_1(p_{\eta|G}, p_{\eta|G'}) = \inf_{Q \in \mathcal{Q}} \int \|\eta_1 - \eta_2\| dQ(\eta_1, \eta_2).$$

Lemma 6. Let $G, G' \subset \Delta^d$ be closed convex subsets such that any $\boldsymbol{\eta} = (\eta_0, \dots, \eta_d) \in G \cup G'$ satisfies $\min_{l=0, \dots, d} \eta_l > c_0$ for some constant $c_0 > 0$. Then

$$K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq \frac{n}{c_0} W_1(p_{\boldsymbol{\eta}|G}, p_{\boldsymbol{\eta}|G'}).$$

Remark. As $n \rightarrow \infty$, the upper bound tends to infinity. This is expected, because the marginal distribution $P_{\mathcal{S}_{[n]}|G}$ should degenerate. Since typically $\text{aff } G \neq \text{aff } G'$, Kullback-Leibler distances between $P_{\mathcal{S}_{[n]}|G}$ and $P_{\mathcal{S}_{[n]}|G'}$ should typically tend to infinity.

Proof. Associating each sample $\mathcal{S}_{[n]} = (X_1, \dots, X_n)$ with a $d + 1$ -dimensional vector $\boldsymbol{\eta}(\mathcal{S}) \in \Delta^d$, where $\boldsymbol{\eta}(\mathcal{S})_i = \frac{1}{n} \sum_{j=1}^n \mathbb{I}(X_j = i)$ for each $i = 0, \dots, d$. The density of $\mathcal{S}_{[n]}$ given G (with respect to the counting measure) takes the form:

$$p_{\mathcal{S}_{[n]}|G}(\mathcal{S}_{[n]}) = \int_G p(\mathcal{S}_{[n]}|\boldsymbol{\eta}) dP(\boldsymbol{\eta}|G) = \int_G \exp\left(n \sum_{i=0}^d \boldsymbol{\eta}(\mathcal{S})_i \log \eta_i\right) dP(\boldsymbol{\eta}|G).$$

Due to the convexity of Kullback-Leibler divergence, by Jensen inequality, for any coupling $Q \in \mathcal{Q}$:

$$\begin{aligned} K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) &= K\left(\int p(\mathcal{S}_{[n]}|\boldsymbol{\eta}_1) dQ(\boldsymbol{\eta}_1, \boldsymbol{\eta}_2), \int p(\mathcal{S}_{[n]}|\boldsymbol{\eta}_2) dQ(\boldsymbol{\eta}_1, \boldsymbol{\eta}_2)\right) \\ &\leq \int K(p(\mathcal{S}_{[n]}|\boldsymbol{\eta}_1), p(\mathcal{S}_{[n]}|\boldsymbol{\eta}_2)) dQ(\boldsymbol{\eta}_1, \boldsymbol{\eta}_2). \end{aligned}$$

It follows that $K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq \inf_Q \int K(p_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_1}, p_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_2}) dQ(\boldsymbol{\eta}_1, \boldsymbol{\eta}_2)$.

Note that $K(P_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_1}, P_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_2}) = \sum_{\mathcal{S}_{[n]}} n(K(\boldsymbol{\eta}(\mathcal{S}), \boldsymbol{\eta}_2) - K(\boldsymbol{\eta}(\mathcal{S}), \boldsymbol{\eta}_1)) p_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_1}$, where the summation is taken over all realizations of $\mathcal{S}_{[n]} \in \{0, \dots, d\}^n$. For any $\boldsymbol{\eta}(\mathcal{S}) \in \Delta^d$, $\boldsymbol{\eta}_1 \in G$ and $\boldsymbol{\eta}_2 \in G'$,

$$\begin{aligned} |K(\boldsymbol{\eta}(\mathcal{S}), \boldsymbol{\eta}_1) - K(\boldsymbol{\eta}(\mathcal{S}), \boldsymbol{\eta}_2)| &= \left| \sum_{i=0}^d \boldsymbol{\eta}(\mathcal{S})_i \log(\eta_{1,i}/\eta_{2,i}) \right| \\ &\leq \sum_i \boldsymbol{\eta}(\mathcal{S})_i |\eta_{1,i} - \eta_{2,i}| / c_0 \\ &\leq \left(\sum_i \boldsymbol{\eta}(\mathcal{S})_i^2 \right)^{1/2} \|\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2\| / c_0 \\ &\leq \|\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2\| / c_0. \end{aligned}$$

Here, the first inequality is due the assumption, the second due to Cauchy-Schwarz. It follows that $K(P_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_1}, P_{\mathcal{S}_{[n]}|\boldsymbol{\eta}_2}) \leq n\|\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2\|/c_0$, so $K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq \frac{n}{c_0} W_1(p_{\boldsymbol{\eta}|G}, p_{\boldsymbol{\eta}|G'})$. \square

Lemma 7. Let $G = \text{conv}(\theta_1, \dots, \theta_k)$ and $G' = \text{conv}(\theta'_1, \dots, \theta'_k)$ (same k). A random variable $\eta \sim P_{\eta|G}$ is parameterized by $\eta = \sum_j \beta_j \eta_j$, while a random variable $\eta \sim P_{\eta|G'}$ is parameterized by $\eta = \sum_j \beta'_j \eta'_j$, where β and β' are both distributed according to a symmetric probability density p_β .

- (a) Assume that both G, G' satisfy property A2. Then, for small $d_{\mathcal{H}}(G, G')$, $W_1(p_{\eta|G}, p_{\eta|G'}) \leq C_0 d_{\mathcal{H}}(G, G')$ for some constant C_0 specified by Lemma 1.
- (b) Assume further that assumptions in Lemma 6 hold, then $K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) \leq \frac{n}{c_0} C_0 d_{\mathcal{H}}(G, G')$.

Remark. In order to obtain an upper bound for $K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'})$ in terms of $d_{\mathcal{H}}(G, G')$, the assumption that p_β is symmetric appears essential. That is, random variables β_1, \dots, β_k are exchangeable under p_β . Without this assumption, it is possible to have $d_{\mathcal{H}}(G, G') = 0$, but $K(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|G'}) > 0$.

Proof. By Lemma 1 under property A2, $d_{\mathcal{M}}(G, G') \leq C_0 d_{\mathcal{H}}(G, G')$ for some constant C_0 . Let $d_{\mathcal{H}}(G, G') \leq \epsilon$ for some small $\epsilon > 0$. Assume without loss of generality that $|\theta_j - \theta'_j| \leq C_0 \epsilon$ for all $j = 1, \dots, k$ (otherwise, simply relabel the subscripts for θ'_j 's).

Let $Q(\eta, \eta')$ be a coupling of $P_{\eta|G}$ and $P_{\eta|G'}$ such that under Q , $\eta = \sum_{j=1}^k \beta_j \theta_j$ and $\eta' = \sum_{j=1}^k \beta_j \theta'_j$, i.e., η and η' share the same β , where β is a random variable with density p_β . This is a valid coupling, since p_β is assumed to be symmetric.

Under distribution Q , $\mathbb{E}\|\eta - \eta'\| \leq \mathbb{E} \sum_{j=1}^k \beta_j \|\theta_j - \theta'_j\| \leq C_0 \epsilon \mathbb{E} \sum_{j=1}^k \beta_j = C_0 \epsilon$. Hence $W_1(P_{\eta|G}, P_{\eta|G'}) \leq C_0 \epsilon$. Part (b) is an immediate consequence. \square

Recall the definition of Kullback-Leibler neighborhood given by Eq. (6). We are now ready to prove the main result of this section:

Theorem 6. Under Assumptions (S1) and (S2), for any G_0 in the support of prior Π , for any $\delta > 0$ and $n > \log(1/\delta)$

$$\Pi(G \in B_K(G_0, \delta)) \geq c(\delta^2/n^3)^{kd},$$

where constant $c = c(c_0, c'_0)$ depends only on c_0, c'_0 .

Proof. We shall invoke a bound of Wong and Shen [1995] (Theorem 5) on the KL divergence. This bound says that if p and q are two densities on a common space such that $\int p^2/q < M$, then for some universal constant $\epsilon_0 > 0$, as long as $h(p, q) \leq \epsilon < \epsilon_0$, there holds: $K(p, q) = O(\epsilon^2 \log(M/\epsilon))$, and $K_2(p, q) := \int p(\log(p/q))^2 = O(\epsilon^2 [\log(M/\epsilon)]^2)$, where the big O constants are universal.

Let $G_0 = \text{conv}(\theta_1^*, \dots, \theta_k^*)$. Consider a random set $G \in \mathcal{G}^k$ represented by $G = \text{conv}(\theta_1, \dots, \theta_k)$, and the event \mathcal{E} that $\|\theta_j - \theta_j^*\| \leq \epsilon$ for all $j = 1, \dots, k$. For the pair of G_0 and G , consider a coupling Q for $P_{\eta|G}$ and $P_{\eta|G_0}$ such that any (η_1, η_2) distributed by Q is parameterized by $\eta_1 = \beta_1 \theta_1 + \dots + \beta_k \theta_k$ and $\eta_2 = \beta_1 \theta_1^* + \dots + \beta_k \theta_k^*$ (that is, under the coupling η_1 and η_2 share the same vector β). Then, under Q , $\mathbb{E}\|\eta_1 -$

$\|\eta_2\| \leq \epsilon$. This entails that $W_1(P_{\eta|G}, P_{\eta|G_0}) \leq \epsilon$. (We note here that the argument appears similar to the one from Lemma 7, but we do not need to assume that p_β be symmetric in this theorem). If G is randomly distributed according to prior Π , under assumption (S2), the probability of event \mathcal{E} is lower bounded by $c'_0 \epsilon^{kd}$. By Lemma 6, $h^2(p_{G_0}, p_G) \leq K(p_{G_0}, p_G)/2 \leq (n/c_0)W_1(P_{\eta|G}, P_{\eta|G_0}) \leq n\epsilon/(2c_0)$. Note that the density ratio $p_{\mathcal{S}_{[n]}|G}/p_{\mathcal{S}_{[n]}|G_0} \leq (1/c_0)^n$, which implies that $\sum_{\mathcal{S}_{[n]}} p_{\mathcal{S}_{[n]}|G_0}^2/p_{\mathcal{S}_{[n]}|G} \leq (1/c_0)^n$. We can apply the upper bound described in the previous paragraph to obtain:

$$K_2(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G}) = O\left(\frac{n\epsilon}{2c_0} \left[\frac{1}{2} \log \frac{2c_0}{n\epsilon} + n \log \frac{1}{c_0} \right]^2\right).$$

Here, the big O constant is universal. If we set $\epsilon = \delta^2/n^3$, then the quantity in the right hand side of the previous display is bounded by $O(\delta^2)$ as long as $n > \log(1/\delta)$. Combining with the probability bound $c'_0 \epsilon^{kd}$ derived above, we obtain the desired result. \square

7 Proofs of main theorems and auxiliary lemmas

Proof of Theorem 1 (Overfitted setting).

Proof. The proof proceeds by verifying conditions of Theorem 4. Let $\epsilon_{m,n} = (\log m/m)^{1/2} + (\log n/m)^{1/2} + (\log n/n)^{1/2}$. Choose the sequence of subsets \mathcal{G}_m simply to be the support of prior Π , so that $\Pi(\mathcal{G}^* \setminus \mathcal{G}_m) = 0$. Note that $\mathcal{G}_m \subset \mathcal{G}^k$. Condition (8) trivially holds. Turning to the entropy conditions, we note that

$$\log D(\epsilon/2, \mathcal{G}_m \cap B_{\mathcal{H}}(G_0, 2\epsilon), d_{\mathcal{H}}) \leq \log N(\epsilon/4, \mathcal{G}_m \cap B_{\mathcal{H}}(G_0, 2\epsilon), d_{\mathcal{H}}) = O(1).$$

By Theorem 5 (a), assumption (S4) and the general inequality that $h \geq V$, we have: $\Psi_{\mathcal{G}_m,n}(\epsilon) \geq [c_1(\epsilon/2)^{p+\alpha} - 6(d+1)e^{-n\epsilon^2/32(d+1)}]^2$, where p is defined as $p = \min(k-1, d)$. So $\Psi_{\mathcal{G}_m,n}(\epsilon) \geq c\epsilon^{2(p+\alpha)}$ as long as $c_1(\epsilon/2)^{p+\alpha} \geq 12(d+1)\exp[-n\epsilon^2/32(d+1)]$. Here, c is a constant depending on c_1, p, d . This is satisfied if ϵ is bounded from below by a large multiple of $\epsilon_{m,n} > (\log n/n)^{1/2}$. Using $\Phi_{\mathcal{G},n}(\delta) := \frac{c_0}{4nC_0} \Psi_{\mathcal{G},n}(\delta)$, it follows that

$$\begin{aligned} & \log D(c_0 \Psi_{\mathcal{G}_m,n}(\epsilon)/(4nC_0), \mathcal{G}_m \cap B_{\mathcal{H}}(G_1, \epsilon/2), d_{\mathcal{H}}) \\ & \leq \log N(c_0 c \epsilon^{2(p+\alpha)}/(4nC_0), \mathcal{G}_m \cap B_{\mathcal{H}}(G_1, \epsilon/2), d_{\mathcal{H}}) \\ & \lesssim \log(n^{kd} \epsilon^{-(2p+2\alpha-1)kd}) \leq m\epsilon^2, \end{aligned}$$

where the last inequality holds since ϵ is bounded from below by a large multiple of $\epsilon_{m,n} > (\log n/m)^{1/2} + (\log m/m)^{1/2}$. Thus, the entropy condition (7)) is established.

To verify condition Eq. (11), we note that for some constant $c > 0$,

$$\begin{aligned} & \exp(2m\epsilon_{m,n}^2) \sum_{j \geq M_m} \exp[-m\Psi_{\mathcal{G}_m,n}(j\epsilon_{m,n})/8] \\ & \leq \exp(2m\epsilon_{m,n}^2) \sum_{j \geq M_m} \exp[-cm(j\epsilon_{m,n})^{2(p+\alpha)}/8] \\ & \lesssim \exp(2m\epsilon_{m,n}^2) \exp[-cm(M_m\epsilon_{m,n})^{2(p+\alpha)}/8], \end{aligned}$$

where the right side of the above display vanishes if $(M_m \epsilon_{m,n})^{p+\alpha}$ is a sufficiently large multiple of $\epsilon_{m,n}$. This holds if we choose $M_m = M \epsilon_{m,n}^{-\frac{p+\alpha-1}{p+\alpha}}$ for a large constant M . Eq. (10) also holds.

It remains to verify Eq. (9). By Theorem 6, as long as $n \gtrsim \log(1/\epsilon_{m,n})$,

$$\log \Pi(G \in B_K(G_0, \epsilon_{m,n})) \geq c(c_0) \log(\epsilon_{m,n}^2/n^3)^{kd} = c(c_0)kd(2 \log \epsilon_{m,n} - 3 \log n).$$

Eq. (9) holds for a sufficiently large constant C because $\epsilon_{m,n} > (\log n/m)^{1/2} + (\log m/m)^{1/2}$, and the constraint that $n > \log m$.

Now, we can apply Theorem 4 to obtain a posterior contraction rate $M_m \epsilon_{m,n} \asymp \epsilon_{m,n}^{1/(p+\alpha)}$. \square

Proof of Theorem 2. The proof proceeds in exactly the same way as Theorem 1, except that part (b) of Theorem 5 is applied instead of part (a). Accordingly p is replaced by 1 in the rate exponent.

Proof of Theorem 3 (Minimax lower bounds). (a) The proof involves the construction of a pair of polytopes in \mathcal{G}^k whose set difference has small volume for a given Hausdorff distance. We consider two separate cases: (i) $k/2 \leq d$ and (ii) $k > 2d$.

If $k/2 \leq d$, consider a $q = \lfloor k/2 \rfloor$ -simplex G_0 that is spanned by $q+1$ vertices in general positions. Take a vertex of G_0 , say θ_0 . Construct G'_0 by chopping G_0 off by an ϵ -cap that is obtained by the convex hull of θ_0 and q other points which lie on the edges adjacent to θ_0 , and of distance ϵ from θ_0 . Clearly, G'_0 has $2q \leq k$ vertices, so both G_0 and G'_0 are in \mathcal{G}^k . We have $d_{\mathcal{H}}(G_0, G'_0) \asymp \epsilon$, and $\text{vol}_q(G_0 \setminus G'_0) \asymp \epsilon^q$. Due to Assumption (S5), $V(p_{\eta|G_0}, p_{\eta|G'_0}) \lesssim \epsilon^{q+\alpha}$. We note here and for the rest of the proof, the multiplying constants in asymptotic inequalities depend only on r, R, δ of properties A1 and A2.

If $k \geq 2d$, consider a d -dimensional polytope G_0 which has $k-d+1$ vertices in general positions. Construct G'_0 in the same way as above (by chopping G_0 off by an ϵ -cap that contains a vertex θ_0 which has d adjacent vertices). Then, G'_0 has $(k-d+1)-1+d = k$ vertices. Thus, both G'_0 and G_0 are in \mathcal{G}^k . We have $d_{\mathcal{H}}(G_0, G'_0) \asymp \epsilon$, and $\text{vol}_d(G_0 \setminus G'_0) \asymp \epsilon^d$. Due to Assumption (S5), $V(p_{\eta|G_0}, p_{\eta|G'_0}) \lesssim \epsilon^{d+\alpha}$.

To combine the two cases, let $q = \min(\lfloor k/2 \rfloor, d)$. We have constructed a pair of $G_0, G'_0 \in \mathcal{G}^k$ such that $d_{\mathcal{H}}(G_0, G'_0) \asymp \epsilon$, and $V(p_{\eta|G_0}, p_{\eta|G'_0}) \lesssim \epsilon^{q+\alpha}$. By Lemma 6, $K(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G'_0}) \lesssim nW_1(p_{\eta|G_0}, p_{\eta|G'_0}) \lesssim nV(p_{\eta|G_0}, p_{\eta|G'_0}) \leq Cn\epsilon^{q+\alpha}$ for some constant $C > 0$ independent of ϵ and n . Note that the second inequality in the above display is due to Theorem 6.15 of Villani [2008].

Applying the method due to Le Cam (cf. Yu [1997], Lemma 1), for any estimator \hat{G} ,

$$\max_{G \in \{G_0, G'_0\}} P_{\mathcal{S}_{[n]}|G_0} d_{\mathcal{H}}(G, \hat{G}) \gtrsim \epsilon \left(1 - \frac{1}{2} V(P_{\mathcal{S}_{[n]}|G_0}^m, P_{\mathcal{S}_{[n]}|G'_0}^m)\right).$$

Here, $P_{\mathcal{S}_{[n]}|G_0}^{[m]}$ denotes the (product) distribution of the m -sample $\mathcal{S}_{[n]}^1, \dots, \mathcal{S}_{[n]}^m$. Thus,

$$\begin{aligned} V^2(P_{\mathcal{S}_{[n]}|G_0}^m, P_{\mathcal{S}_{[n]}|G'_0}^m) &\leq h^2(P_{\mathcal{S}_{[n]}|G_0}^m, P_{\mathcal{S}_{[n]}|G'_0}^m) \\ &= 1 - \int [P_{\mathcal{S}_{[n]}|G_0}^m P_{\mathcal{S}_{[n]}|G'_0}^m]^{1/2} \\ &= 1 - [1 - h^2(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G'_0})]^m \\ &\leq 1 - (1 - Cn\epsilon^{q+\alpha})^m. \end{aligned}$$

The last inequality is due to $h^2(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G'_0}) \leq K(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G'_0}) \leq Cn\epsilon^{q+\alpha}$. Thus,

$$\max_{G \in \{G_0, G'_0\}} P_{\mathcal{S}_{[n]}|G_0} d_{\mathcal{H}}(G, \hat{G}) \gtrsim \epsilon(1 - \frac{1}{2}[1 - (1 - Cn\epsilon^{q+\alpha})^m]^{1/2}).$$

Letting $\epsilon^{q+\alpha} = \frac{1}{Cmn}$, the right side of the previous display is bounded from below by $\epsilon(1 - \frac{1}{2}(1 - 1/2)^{1/2})$.

(b) We employ the same construction of G_0 and G'_0 as in part (a). Using the argument used in the proof of Lemma 5 $K(p_{\mathcal{S}_{[n]}|G'_0}, p_{\mathcal{S}_{[n]}|G_0}) = \int \log[\text{vol}_q G_0 / \text{vol}_q G'_0] dP_{\boldsymbol{\eta}|G_0} \leq \int \log(1 + C\epsilon^q) P_{\boldsymbol{\eta}|G_0} \lesssim \epsilon^q$. So, $h^2(p_{\mathcal{S}_{[n]}|G_0}, p_{\mathcal{S}_{[n]}|G'_0}) \leq K(p_{\mathcal{S}_{[n]}|G'_0}, p_{\mathcal{S}_{[n]}|G_0}) \lesssim \epsilon^q$. Then, the proof proceeds as in part (a).

(c) Let G'_0 be a polytope such that $|\text{extr } G'_0| = |\text{extr } G_0| = k$ and $d_{\mathcal{H}}(G'_0, G_0) = \epsilon$. By Lemma 2, $\text{vol}_p(G_0 \triangle G'_0) = O(\epsilon)$, where $p = (k-1) \wedge d$. The proof proceeds as in part (a) to obtain $(1/mn)^{1/(1+\alpha)}$ rate for the lower bound under assumption (S5). Under assumption (S5'), as in part (b), the dependence on n can be removed to obtain $1/m$ rate.

Proof of α -regularity of the Dirichlet-induced densities in Lemma 4.

Proof. First, consider the case $k \leq d+1$. For $\boldsymbol{\eta}^* \in G$, write $\boldsymbol{\eta}^* = \beta_1^* \boldsymbol{\theta}_1 + \dots + \beta_k^* \boldsymbol{\theta}_k$. For $\boldsymbol{\beta} \in \Delta^{k-1}$ such that $|\beta_i - \beta_i^*| \leq \epsilon/k$ for all $i = 1, \dots, k-1$, we have $\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| = \|\sum_{i=1}^k (\beta_i - \beta_i^*) \boldsymbol{\theta}_i\| \leq \sum_{i=1}^k |\beta_i - \beta_i^*| \leq 2 \sum_{i=1}^{k-1} |\beta_i - \beta_i^*| \leq 2\epsilon$. Here, we used the fact that $\|\boldsymbol{\theta}_i\| \leq 1$ for any $\boldsymbol{\theta}_i \in \Delta^d$. Without loss of generality, assume that $\beta_k^* \geq 1/k$. Then, for any $\epsilon < 1/k$

$$\begin{aligned} P_{\boldsymbol{\eta}|G}(\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| \leq 2\epsilon) &\geq P_{\boldsymbol{\beta}}(|\beta_i - \beta_i^*| \leq \epsilon/k; i = 1, \dots, k-1) \\ &= \frac{\Gamma(\sum \gamma_i)}{\prod_i \Gamma(\gamma_i)} \int_{\beta_i \in [0,1]; |\beta_i - \beta_i^*| \leq \epsilon/k; i=1, \dots, k-1} \prod_{i=1}^{k-1} \beta_i^{\gamma_i-1} (1 - \sum_{i=1}^{k-1} \beta_i)^{\gamma_k-1} d\beta_1 \dots d\beta_{k-1} \\ &\geq \frac{\Gamma(\sum \gamma_i)}{\prod_i \Gamma(\gamma_i)} \prod_{i=1}^{k-1} \int_{\max(\gamma_i^* - \epsilon/k, 0)}^{\min(\gamma_i^* + \epsilon/k, 1)} \beta_i^{\gamma_i-1} d\beta_i \geq \frac{\Gamma(\sum \gamma_i)}{\prod_i \Gamma(\gamma_i)} (\epsilon/k)^{k-1}. \end{aligned}$$

Both the second and the third inequality in the previous display exploits the fact that since $\gamma_i \leq 1$, $x^{\gamma_i-1} \geq 1$ for any $x \leq 1$.

Now, consider the case $k > d+1$. The proof in the previous case applies, but we can achieve a better lower bound because the intrinsic dimensionality of G is d , not $k-1$. Since

$\boldsymbol{\eta}^* \in \text{conv}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k) \subset \Delta^d$, by Carathéodory's theorem, $\boldsymbol{\eta}^*$ is the convex combination of $d + 1$ or fewer extreme points among $\boldsymbol{\theta}_i$'s. Without loss of generality, let $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{d+1}$ be such points, and write $\boldsymbol{\eta}^* = \beta_1^* \boldsymbol{\theta}_1 + \dots + \beta_{d+1}^* \boldsymbol{\theta}_{d+1}$. Consider $\boldsymbol{\eta} = \beta_1 \boldsymbol{\theta}_1 + \dots + \beta_k \boldsymbol{\theta}_k$, where $\|\beta_i - \beta_i^*\| \leq \epsilon/k$, for $i = 1, \dots, d$, while $0 \leq \beta_i \leq \epsilon/k$ for $i = d+2, \dots, k$. Then, $\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| \leq 2\epsilon$. This implies that

$$\begin{aligned} P_{\boldsymbol{\eta}|G}(\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| \leq 2\epsilon) &\geq P_{\boldsymbol{\beta}}(|\beta_i - \beta_i^*| \leq \epsilon/k, i = 1, \dots, d+1; |\beta_j| \leq \epsilon/k, j > d+1) \\ &\geq \frac{\Gamma(\sum \gamma_i)}{\prod_i \Gamma(\gamma_i)} \prod_{i=1}^d \int_{\max(\gamma_i^* - \epsilon/k, 0)}^{\min(\gamma_i^* + \epsilon/k, 1)} \beta_i^{\gamma_i - 1} d\beta_i \prod_{i=d+2}^k \int_0^{\epsilon/k} \beta_i^{\gamma_i - 1} d\beta_i \\ &\geq \frac{\Gamma(\sum \gamma_i)}{\prod_i \Gamma(\gamma_i)} (\epsilon/k)^{d + \sum_{i=d+2}^k \gamma_i} / \prod_{i=d+2}^n \gamma_i \gtrsim \epsilon^{d + \sum_{i=1}^k \gamma_i}. \end{aligned}$$

This concludes the proof. \square

8 Appendix A: Proofs of geometric lemmas

Proof of Lemma 1

Proof. (a) Let $G = \text{conv}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)$ and $G' = \text{conv}(\boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_{k'})$. This part of the lemma is immediate from the definition by noting that for any $x \in G$, $d(x, G') \leq \min_j \|\boldsymbol{\theta}_j - \boldsymbol{\theta}'_j\|$, while the maximum of $d(x, G')$ is attained at some extreme point of G .

(b) Let $d_{\mathcal{H}}(G, G') = \epsilon$ for some small $\epsilon > 0$. Take an extreme point of G , say $\boldsymbol{\theta}_1$. Due to A2, there is a ray emanating from $\boldsymbol{\theta}_1$ that intersects with the interior of G and the angles formed by the ray and all (exposed) edges incident to $\boldsymbol{\theta}_1$ are bounded from above by $\pi/2 - \delta$. Let x be the intersection between the ray and the boundary of $B_p(\boldsymbol{\theta}_1, \epsilon)$.

Let H be a $p-1$ -dimensional hyperplane in \mathbb{R}^p that touches (intersects with) $B_p(\boldsymbol{\theta}_1, \epsilon)$ at only x . Define $C(x)$, resp. $C_\epsilon(x)$, to be the p -dimensional caps obtained by the intersection between G , resp. G_ϵ , with the half-space which contains $\boldsymbol{\theta}_1$ and which is supported by H . For any x' that lies in the intersection of H and a line segment $[\boldsymbol{\theta}_1, \boldsymbol{\theta}_i]$, where $\boldsymbol{\theta}_i$ is another vertex of G , the line segment $[x, x'] \in H$ and $\|x - x'\| \leq \epsilon \cot \delta$. Suppose that the ray emanating from x through x' intersects with $\text{bd } G_\epsilon$ at x'' . Then, $\|x' - x''\| \leq \epsilon / \sin \delta$, which implies that $\|x - x''\| \leq \epsilon(\cot \delta + 1/\sin \delta)$ by triangle inequality. This entails that $\text{Diam } C_\epsilon(x) \leq C\epsilon$, where $C = (1 + (\cot \delta + 1/\sin \delta)^2)^{1/2}$.

Now, $d_{\mathcal{H}}(G, G') = \epsilon$ implies that $G' \cap B_p(\boldsymbol{\theta}_1, \epsilon) \neq \emptyset$. There is an extreme point of G' in the half-space which contains $B(\boldsymbol{\theta}_1, \epsilon)$ and is supported by H . But $G' \subset G_\epsilon$, so there is an extreme point of G' in $C_\epsilon(x)$. Hence, there is $\boldsymbol{\theta}'_j \in G'$ such that $\|\boldsymbol{\theta}'_j - \boldsymbol{\theta}_1\| \leq \text{Diam}(C_\epsilon(x)) \leq C\epsilon$. Repeat this argument for all other extreme points of G to conclude that $d_{\mathcal{M}}(G, G') \leq C\epsilon$. \square

Proof of Lemma 2

Proof. (a) Let $d_{\mathcal{H}}(G, G') = \epsilon$. There exists either a point $x \in \text{bd } G$ such that $G' \cap B_p(x, \epsilon/2) = \emptyset$, or a point $x' \in \text{bd } G'$ such that $G \cap B_p(x', \epsilon/2) = \emptyset$. Without loss of generality, assume the former. Thus, $\text{vol}_p G \triangle G' \geq \text{vol}_p B_p(x, \epsilon/2) \cap G$. Consider the convex cone emanating from x that circumscribes the p -dimensional spherical ball $B_p(\theta_c, r)$ (whose existence is given by Condition A1). Since $\|x - \theta_c\| \leq R$, the angle between the line segment $[x, \theta_c]$ and the cone's rays is bounded from below by $\sin \varphi \geq r/R$. So, $\text{vol}_p B_d(x, \epsilon/2) \cap G \geq c_1 \epsilon^p$, where c_1 depends only on r, R, p .

(b) Let $d_{\mathcal{H}}(G, G') = \epsilon$. Then $G' \subset G_\epsilon$ and $G \subset G'_\epsilon$. Take any point $x \in \text{bd } G$, let x' be the intersection between $\text{bd } G_\epsilon$ and the ray emanating from θ_c and passing through x . Let H_1 be a $p-1$ dimensional supporting hyperplane for G at x . There is also a supporting hyperplane H_2 of G' that is parallel to H_1 and of at most ϵ distance away from H_1 . Since $\|\theta_c - x\| \leq R$, while the distance from θ_c to H_1 is lower bounded by r , the angle φ between vector $\theta_c - x$ and the vector normal to H_1 satisfies $\cos \varphi \geq r/R$. This implies that $\|x' - x\| \leq \epsilon / \cos \varphi \leq \epsilon R / r$, so $\|x' - \theta_c\| / \|x - \theta_c\| \leq 1 + \epsilon R / r^2$. In other words, $G_\epsilon - \theta_c \subset (1 + \epsilon R / r^2)(G - \theta_c)$. So, $\text{vol}_p G' \setminus G \leq \text{vol}_p G_\epsilon \setminus G \leq [(1 + \epsilon R / r^2)^p - 1] \text{vol}_p G \leq C_1 \epsilon$, where C_1 depends only on r, R, p . We obtain a similar bound for $\text{vol}_p G \setminus G'$, which concludes the proof. \square

Proof of Lemma 3

Proof. We provide a proof for case (a). Let $G = \text{conv}(\theta_1, \dots, \theta_k)$ and $G' = \text{conv}(\theta'_1, \dots, \theta'_k)$, where G is fixed but G' is allowed to vary. Since G is fixed, it satisfies A1 and A2 for some constants r, R and δ (depending on G). Moreover, there is some $\epsilon_0 = \epsilon_0(G)$ depending only on G such that as soon as $d_{\mathcal{H}}(G, G') \leq \epsilon_0$, G' also satisfies A1 and A2 for constants $\delta' = \delta/2, r' = r/2, R' = 2R$.

Suppose that $d_{\mathcal{H}}(G, G') = \epsilon$ such that $\epsilon < \epsilon_0$. By Lemma 1 (b) for each vertex of G , say θ_i , there is a vertex of G' , say θ'_i , such that $\theta'_i \in B_p(\theta_i, C_0 \epsilon)$ with $C_0 = C_0(G)$ depending only on δ . Moreover, there is at least one vertex of G , say θ_1 , for which $\|\theta'_1 - \theta_1\| \geq \epsilon$.

There are only three possible general positions for θ'_1 relatively to G . Either

- (i) $\theta'_1 \in G$, or
- (ii) $\theta'_1 \in 2\theta_1 - G$, or
- (iii) θ'_1 lies in a cone formed by all half-spaces supported by the $p-1$ dimensional faces adjacent to θ_1 . Among these there is one half-space that contains G , and one that does not contain G .

If (i) is true, by property A1, G has at least one face $S \supset \theta_1$ such that the distance from θ'_1 to the hyperplane that provides support for S is bounded from below by $\epsilon r / R$. Let $B \subset S$ be a homothetic transformation of S with respect to center θ_1 that maps $x \in S$ to $\tilde{x} \in B$ such that the ratio $\eta := \|\theta_1 - \tilde{x}\| / \|\theta_1 - x\|$ satisfies $1 - \eta = 2C_0 \epsilon / \min_{i \neq j} \|\theta_i - \theta_j\| \in (0, 1/2)$. This is possible as soon as $\epsilon < \min_{i \neq j} \|\theta_i - \theta_j\| / 4C_0$. Then, for any $\theta_j \in S, j \neq 1$, under

this transformation $\theta_j \mapsto \tilde{\theta}_j \in S$ for which $\|\tilde{\theta}_j - \theta_j\| = (1 - \eta)\|\theta_1 - \theta_j\| \geq 2C_0\epsilon$. Since $\|\theta'_j - \theta_j\| \leq C_0\epsilon$, the construction of B implies that $\theta'_j \notin B$. As a result, $B \cap G' = \emptyset$. Moreover, $\text{vol}_{p-1} B = \eta^{p-1} \text{vol}_{p-1} S \geq (1/2)^{p-1} \text{vol}_{p-1} S \geq c_0(G)$, a constant depending only on G . Let Q be a p -pyramid which has apex θ'_1 and base B . It follows that $\text{relint } Q \cap \text{relint } G' = \emptyset$, which implies that $\text{relint } Q \subset G \setminus G'$. (relint stands for the relative interior of a set). Hence, $\text{vol}_p G \setminus G' \geq \text{vol}_p Q \geq \frac{1}{p}\epsilon r/R \text{vol}_{p-1} B \geq \frac{1}{p}\epsilon c_0(G)r/R$.

If (ii) is true, the same argument can be applied to show that $\text{vol}_p(G' \setminus G) = \Omega(\epsilon)$. If (iii) is true, a similar argument continues to apply: we obtain a lower bound for either $\text{vol}_p G' \setminus G$ or $\text{vol}_p G \setminus G'$. G has a face (supported by a hyperplane, say, H) such that the distance from θ'_1 to H is $\Omega(\epsilon)$. If the half-space supported by H that contains θ'_1 but does not contain G , then $\text{vol}_p G' \setminus G = \Omega(\epsilon)$. If, on the other hand, the associated half-space does contain G , then $\text{vol}_p G \setminus G' = \Omega(\epsilon)$. The proof for case (b) is similar and is omitted. \square

9 Appendix B: Proof of abstract posterior contraction theorem

A key ingredient in the general analysis of convergence of posterior distributions is through establishing the existence of tests for subsets of parameters of interest. A test $\varphi_{m,n}$ is a measurable indicator function of the $m \times n$ -sample $\mathcal{S}_{[n]}^{[m]} = (\mathcal{S}_{[n]}^1, \dots, \mathcal{S}_{[n]}^m)$ from an admixture model. For a fixed pair of convex polytopes $G_0, G_1 \in \mathcal{G}$, where \mathcal{G} is a given subset of Δ^d , consider tests for discriminating G_0 against a closed Hausdorff ball centered at G_1 . The following two lemmas on the existence of tests highlight the fundamental role of the Hellinger information:

Lemma 8. *Fix a pair of $(G_0, G_1) \in (\mathcal{G}^* \times \mathcal{G})$ and let $\delta = d_{\mathcal{H}}(G_0, G_1)$. Then, there exist tests $\{\varphi_{m,n}\}$ that have the following properties:*

$$P_{\mathcal{S}_{[n]}|G_0}^m \varphi_{m,n} \leq D \exp[-m\Psi_{\mathcal{G},n}(\delta)/8] \quad (16)$$

$$\sup_{G \in \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)} P_{\mathcal{S}_{[n]}|G}^m (1 - \varphi_{m,n}) \leq \exp[-m\Psi_{\mathcal{G},n}(\delta)/8]. \quad (17)$$

Here, $D := D\left(\Phi_{\mathcal{G},n}(\delta), \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2), d_{\mathcal{H}}\right)$, i.e., the maximal number of elements in $\mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)$ that are mutually separated by at least $\Phi_{\mathcal{G},n}(\delta)$ in Hausdorff metric $d_{\mathcal{H}}$.

Proof. We begin the proof by noting that a direct application of standard results on existence of tests (cf. Cam [1986], Chapter 4) is not possible, due to the lack of convexity of the space of densities of $\mathcal{S}_{[n]}$ as G varies in some subset $\mathcal{G} \subset \mathcal{G}^*$, even if \mathcal{G} is convex. This difficulty is overcome by appealing to a packing argument.

Consider a maximal $\Phi_{\mathcal{G},n}(\delta)$ -packing in $d_{\mathcal{H}}$ metric for the set $\mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)$. This yields a set of $D = D(\Phi_{\mathcal{G},n}(\delta), \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2), d_{\mathcal{H}})$ elements $\tilde{G}_1, \dots, \tilde{G}_D \in \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)$.

Next, we note the following fact: for any $t = 1, \dots, D$, if $G \in \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)$ and $d_{\mathcal{H}}(G, \tilde{G}_t) \leq \Phi_{\mathcal{G},n}(\delta)$, then by the definition of Φ , $h^2(p_{\mathcal{S}_{[n]}|G}, p_{\mathcal{S}_{[n]}|\tilde{G}_t}) \leq \frac{1}{4}\Psi_{\mathcal{G},n}(\delta)$. By

the definition of Hellinger information, $h^2(p_{S_{[n]}|G_0}, p_{S_{[n]}|\tilde{G}_t}) \geq \Psi_{\mathcal{G},n}(\delta)$. Thus, by triangle inequality, $h(p_{S_{[n]}|G_0}, p_{S_{[n]}|G}) \geq \frac{1}{2}\Psi_{\mathcal{G},n}(\delta)^{1/2}$.

For each pair of G_0, \tilde{G}_t there exist tests $\omega_n^{(t)}$ of $p_{S_{[n]}|G_0}$ versus the Hellinger ball $\mathcal{P}_2(t) := \{p_{S_{[n]}|G} | G \in \mathcal{G}^*; h(p_{S_{[n]}|G}, p_{S_{[n]}|\tilde{G}_t}) \leq \frac{1}{2}h(p_{S_{[n]}|G_0}, p_{S_{[n]}|\tilde{G}_t})\}$ such that,

$$\begin{aligned} P_{S_{[n]}|G_0}^m \omega_{m,n}^{(t)} &\leq \exp[-mh^2(p_{S_{[n]}|G_0}, p_{S_{[n]}|\tilde{G}_t})/8], \\ \sup_{P_2 \in \mathcal{P}_2(t)} P_2^m (1 - \omega_{m,n}^{(t)}) &\leq \exp[-mh^2(p_{S_{[n]}|G_0}, p_{S_{[n]}|\tilde{G}_t})/8]. \end{aligned}$$

Consider the test $\varphi_{m,n} = \max_{1 \leq t \leq D} \omega_{m,n}^{(t)}$, then

$$\begin{aligned} P_{S_{[n]}|G_0}^m \varphi_{m,n} &\leq D \times \exp[-m\Psi_{\mathcal{G},n}(\delta)/8], \\ \sup_{G \in \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)} P_{S_{[n]}|G}^m (1 - \varphi_{m,n}) &\leq \exp[-m\Psi_{\mathcal{G},n}(\delta)/8]. \end{aligned}$$

The first inequality is due to $\varphi_{m,n} \leq \sum_{t=1}^D \omega_{m,n}^{(t)}$, and the second is due to the fact that for any $G \in \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \delta/2)$ there is some $d = 1, \dots, D$ such that $d_{\mathcal{H}}(G, \tilde{G}_t) \leq \Phi_{\mathcal{G},n}(\delta)$, so that $p_{S_{[n]}|G} \in \mathcal{P}_2(t)$. \square

Next, the existence of tests can be shown for discriminating G_0 against the complement of a closed Hausdorff ball:

Lemma 9. Suppose that \mathcal{G} satisfies condition C. Fix $G_0 \in \mathcal{G}^k$. Suppose that for some non-increasing function $D(\epsilon)$, some $\epsilon_{m,n} \geq 0$ and every $\epsilon > \epsilon_{m,n}$,

$$\begin{aligned} &\sup_{G_1 \in \mathcal{G}} D(\Phi_{\mathcal{G},n}(\epsilon), \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, \epsilon/2), d_{\mathcal{H}}) \\ &\times D(\epsilon/2, \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_0, 2\epsilon) \setminus B_W(G_0, \epsilon), d_{\mathcal{H}}) \leq D(\epsilon). \end{aligned} \quad (18)$$

Then, for every $\epsilon > \epsilon_{m,n}$, and any $t_0 \in \mathbb{N}$, there exist tests $\varphi_{m,n}$ (depending on $\epsilon > 0$) such that

$$P_{G_0} \varphi_{m,n} \leq D(\epsilon) \sum_{t=t_0}^{\lceil \text{Diam}(\mathcal{G})/\epsilon \rceil} \exp[-m\Psi_{\mathcal{G},n}(t\epsilon)/8] \quad (19)$$

$$\sup_{G \in \mathcal{G}: d_{\mathcal{H}}(G_0, G) > t_0\epsilon} P_G (1 - \varphi_{m,n}) \leq \exp[-m\Psi_{\mathcal{G},n}(t_0\epsilon)/8]. \quad (20)$$

Proof. The proof consists of a standard peeling device (e.g., Ghosal et al. [2000]) and a packing argument as in the previous proof. For a given $t \in \mathbb{N}$ choose a maximal $t\epsilon/2$ -packing for set $S_t = \{G : t\epsilon < d_{\mathcal{H}}(G_0, G) \leq (t+1)\epsilon\}$. This yields a set S'_t of at most $D(t\epsilon/2, S_t, d_{\mathcal{H}})$ points. Moreover, every $G \in S_t$ is within distance $t\epsilon/2$ of at least one of the points in S'_t . For every such point $G_1 \in S'_t$, there exists a test $\omega_{m,n}$ satisfying Eqs. (16) and (17), where δ is taken to be $\delta = t\epsilon$. Take $\varphi_{m,n}$ to be the maximum of all tests

attached this way to some point $G_1 \in S'_t$ for some $t \geq t_0$. Note that $G \in \mathcal{G} \subset \Delta^d$, so $t \leq \lceil \text{Diam}(\mathcal{G})/\epsilon \rceil$. Then, by union bound, and the condition that $D(\epsilon)$ is non-increasing,

$$\begin{aligned} P_{\mathcal{S}_{[n]}|G_0}^m \varphi_{m,n} &\leq \sum_{t=t_0}^{\lceil \frac{\text{Diam}(\mathcal{G})}{\epsilon} \rceil} \sum_{G_1 \in S'_t} D\left(\Phi_{\mathcal{G},n}(t\epsilon), \mathcal{G} \cap B_{d_{\mathcal{H}}}(G_1, t\epsilon/2), d_{\mathcal{H}}\right) \times \\ &\quad \exp[-m\Psi_{\mathcal{G},n}(t\epsilon)/8] \\ &\leq D(\epsilon) \sum_{t \geq t_0} \exp[-m\Psi_{\mathcal{G},n}(t\epsilon)/8], \end{aligned}$$

and

$$\begin{aligned} \sup_{G \in \cup_{u \geq t_0} S_u} P_{\mathcal{S}_{[n]}|G}^m (1 - \varphi_n) &\leq \sup_{u \geq t_0} \exp[-m\Psi_{\mathcal{G},n}(u\epsilon)/8] \\ &\leq \exp[-m\Psi_{\mathcal{G},n}(t_0\epsilon)/8], \end{aligned}$$

where the last inequality is due the monotonicity of $\Psi_{\mathcal{G},n}(\cdot)$. \square

Proof of the abstract posterior contraction theorem (Theorem 4)

Proof. In this proof, to simplify notations denote $P_G := P_{\mathcal{S}_{[n]}|G}$ and so on. By a result of Ghosal et al [Ghosal et al., 2000] (Lemma 8.1, pg. 524), for every $\epsilon > 0, C > 0$ and every probability measure Π_0 supported on the set $B_K(G_0, \epsilon)$ defined by Eq. (6), we have,

$$P_{G_0} \left(\int \prod_{i=1}^m \frac{p_G(\mathcal{S}_{[n]}^i)}{p_{G_0}(\mathcal{S}_{[n]}^i)} d\Pi_0(G) \leq \exp(-(1+C)m\epsilon^2) \right) \leq \frac{1}{C^2 m \epsilon^2}.$$

This entails that, by fixing $C = 1$, there is an event A_m with $P_{G_0}^m$ -probability at least $1 - (m\epsilon_{m,n}^2)^{-1}$, for which there holds:

$$\int \prod_{i=1}^n p_G(\mathcal{S}_{[n]}^i)/p_{G_0}(\mathcal{S}_{[n]}^i) d\Pi(G) \geq \exp(-2m\epsilon_{m,n}^2) \Pi(B_K(G_0, \epsilon_{m,n})). \quad (21)$$

Let $\mathcal{O}_m = \{G \in \mathcal{G}^* : d_{\mathcal{H}}(G_0, G) \geq M_m \epsilon_{m,n}\}$. Due to Eq.(7), the condition specified by Lemma 9 is satisfied by setting $D(\epsilon) = \exp(m\epsilon_{m,n}^2)$ (constant in ϵ). Thus there exist tests $\varphi_{m,n}$ for which Eq. (19) and (20) hold. Then,

$$\begin{aligned} &P_{G_0} \Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]}) \\ &= P_{G_0} [\varphi_{m,n} \Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]})] + P_{G_0} [(1 - \varphi_{m,n}) \Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]})] \\ &\leq P_{G_0} [\varphi_{m,n} \Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]})] + P_{G_0} \mathbb{I}(A_m^c) \\ &\quad + P_{G_0} \left[(1 - \varphi_{m,n}) \Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]}) \mathbb{I}(A_m) \right]. \end{aligned}$$

Applying Lemma 9, the first term in the preceeding display is bounded above by

$$P_{G_0}\varphi_{m,n} \leq D(\epsilon_{m,n}) \sum_{j \geq M_m} \exp[-m\Psi_{\mathcal{G}_{m,n}}(j\epsilon_{m,n})/8] \rightarrow 0,$$

thanks to Eq. (11). The second term in the above display is bounded by $(m\epsilon_{m,n}^2)^{-1}$ by the definition of A_m , so this term vanishes. It remains to show that third term in the display also vanishes as $m \rightarrow \infty$. By Bayes' rule,

$$\Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]}) = \frac{\int_{\mathcal{O}_m} \prod_{i=1}^m p_G(\mathcal{S}_{[n]}^i) / p_{G_0}(\mathcal{S}_{[n]}^i) d\Pi(G)}{\int \prod_{i=1}^m p_G(\mathcal{S}_{[n]}^i) / p_{G_0}(\mathcal{S}_{[n]}^i) d\Pi(G)},$$

and then obtain a lower bound for the denominator by Eq. (21). For the nominator, by Fubini's theorem:

$$\begin{aligned} & P_{G_0} \int_{\mathcal{O}_m \cap \mathcal{G}_m} (1 - \varphi_{m,n}) \prod_{i=1}^m p_G(\mathcal{S}_{[n]}^i) / p_{G_0}(\mathcal{S}_{[n]}^i) d\Pi(G) \\ &= \int_{\mathcal{O}_m \cap \mathcal{G}_m} P_G(1 - \varphi_{m,n}) d\Pi(G) \leq \exp[-m\Psi_{\mathcal{G}_{m,n}}(M_m\epsilon_{m,n})/8], \end{aligned} \quad (22)$$

where the last inequality is due to Eq. (20). In addition, by (8),

$$\begin{aligned} & P_{G_0} \int_{\mathcal{O}_m \setminus \mathcal{G}_m} (1 - \varphi_{m,n}) \prod_{i=1}^m p_G(\mathcal{S}_{[n]}^i) / p_{G_0}(\mathcal{S}_{[n]}^i) d\Pi(G) \\ &= \int_{\mathcal{O}_m \setminus \mathcal{G}_m} P_G(1 - \varphi_{m,n}) d\Pi(G) \leq \Pi(\mathcal{G}^* \setminus \mathcal{G}_m). \end{aligned} \quad (23)$$

Now, combining bounds (22) and (23) with condition (10), we obtain:

$$\begin{aligned} & P_{G_0}(1 - \varphi_{m,n})\Pi(G \in \mathcal{O}_m | \mathcal{S}_{[n]}^{[m]})\mathbb{I}(A_m) \\ &\leq \frac{\Pi(\mathcal{G}^* \setminus \mathcal{G}_m) + \exp[-m\Psi_{G_0,n}(\mathcal{G}_m, M_m\epsilon_{m,n})/8]}{\exp(-2m\epsilon_{m,n}^2)\Pi(B_K(G_0, \epsilon_{m,n}))} \end{aligned}$$

The upper bound in the preceeding display converges to 0 by Eq. (11), thereby concluding the proof. \square

References

- A. Anandkumar, D. Foster, D. Hsu, S. Kakade, and Y. K. Liu. Two SVDs suffice: Spectral decompositions for probabilistic topic modeling and latent Dirichlet allocation. *arXiv:1204.6703*, 2012.

- S. Arora, R. Ge, and A. Moitra. Learning topic models – going beyond SVD. *arXiv:1204.1956*, 2012.
- A. Barron, M. Schervish, and L. Wasserman. The consistency of posterior distributions in nonparametric problems. *Ann. Statist.*, 27:536–561, 1999.
- D.M. Blei, A.Y. Ng, and M.I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3: 993–1022, 2003.
- L. Le Cam. *Asymptotic methods in statistical decision theory*. Springer-Verlag, 1986.
- J. Chen. Optimal rate of convergence for finite mixture models. *Annals of Statistics*, 23(1): 221–233, 1995.
- L. Evans and R. Gariepy. *Measure theory and fine properties of functions*. CRC Press, 1992.
- S. Ghosal and A. van der Vaart. Convergence rates of posterior distributions for noniid observations. *Ann. Statist.*, 35(1):192–223, 2007.
- S. Ghosal, J. K. Ghosh, and A. van der Vaart. Convergence rates of posterior distributions. *Ann. Statist.*, 28(2):500–531, 2000.
- J. K. Ghosh and R. V. Ramamoorthi. *Bayesian nonparametrics*. Springer, 2002.
- H. Ishwaran, L. James, and J. Sun. Bayesian model selection in finite mixtures by marginal density decompositions. *Journal of American Statistical Association*, 96(456):1316–1332, 2001.
- X. Nguyen. Convergence of latent mixing measures in finite and infinite mixture models. *Annals of Statistics*, to appear, 2012.
- J. Pritchard, M. Stephens, and P. Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155:945–959, 2000.
- J. Rousseau and K. Mengersen. Asymptotic behaviour of the posterior distribution in overfitted mixture models. *Journal of the Royal Statistical Society: Series B*, 73(5):689–710, 2011.
- R. Schneider. *Convex bodies: Brunn-Minkowsky theory*. Cambridge University Press, 1993.
- X. Shen and L. Wasserman. Rates of convergence of posterior distributions. *Ann. Statist.*, 29:687–714, 2001.
- W. Toussile and E. Gassiat. Model based clustering using multilocus data with loci selection. *Advances in Data Analysis and Classification*, 3:109–134, 2009.
- S. van de Geer. *Empirical processes in M-estimation*. Cambridge University Press, 2000.

- Cédric Villani. *Optimal transport: Old and New*. Springer, 2008.
- S. Walker. New approaches to bayesian consistency. *Ann. Statist.*, 32(5):2028–2043, 2004.
- S. Walker, A. Lijoi, and I. Prunster. On rates of convergence for posterior distributions in infinite-dimensional models. *Ann. Statist.*, 35(2):738–746, 2007.
- W. H. Wong and X. Shen. Probability inequalities for likelihood ratios and convergences of sieves mles. *Ann. Statist.*, 23:339–362, 1995.
- B. Yu. Assouad, Fano, and Le Cam. *Festschrift for Lucien Le Cam*, pages 423–435, 1997.